

Orthogonal Recursive Bisection on the GPU for
Accelerated Load Balancing in Large N-Body
Simulations

-

Bachelor Thesis

Andrin Rehmann

August 4, 2022

Contents

1	Introduction	4
2	Background	5
2.1	Fast Multipole Expansion	5
2.1.1	Advantage of Cubic Bounding Volumes	5
2.2	ORB vs. K-d Trees	6
2.3	Rise of the GPU	6
2.4	Related Work	7
3	Orthogonal Recursive Bisection (ORB)	8
3.1	Target Data and Notation	8
3.2	Memory and Workload Balancing	9
3.3	ORB	10
3.3.1	Algorithm	11
3.4	By Example	12
3.5	Root finding	15
3.5.1	Bisection Method	16
3.5.2	Edge Cases	17
3.5.3	Runtime Analysis	18
3.6	Partition Algorithm	18
4	Theoretical Analysis	20
4.1	General Memory Model	20
4.2	Supercomputers	21
4.3	Roofline Performance Model	22
4.3.1	Estimating Flops	22
4.3.2	Estimating Arithmetic Intensity	23
4.3.3	The Plots	24
4.3.4	Empirical Verification	27
4.4	Runtime Estimates	27
4.4.1	CPU Version	28
4.4.2	GPU Counting	28
4.4.3	GPU Counting and Partitioning	29
4.4.4	Plugin Values	29
4.5	Conclusion	30
5	CPU Implementation	31
5.1	MDL and PKDGrav	31
5.2	Particles	31
5.3	Cell	32
5.3.1	Heap Conditions	32
5.3.2	Class	33
5.4	Mapping Cells to Particles	34
5.5	Services	34

5.5.1	Init and Finalize Service	34
5.5.2	Count Left Service	34
5.5.3	Count	35
5.5.4	Partition	36
5.5.5	Make Axis	37
5.6	Parallel Schedule	37
6	GPU Implementation	40
6.1	Relevant CUDA Concepts	40
6.1.1	Warps	40
6.1.2	Memory	41
6.1.3	Memory Access Patterns	41
6.1.4	Asynchronous Operations	42
6.1.5	Synchronization	42
6.2	Streams	43
6.3	Memory Management	43
6.4	GPU Accelerated Count Left	43
6.4.1	Schedule	43
6.4.2	Service	45
6.4.3	Kernel Code	46
6.5	Improved GPU Accelerated Count Left	48
6.5.1	Schedule	48
6.5.2	Kernel Code	49
6.6	GPU Accelerated Partitioning	52
6.6.1	Schedule	53
6.6.2	Partitioning a Single Cell	54
6.6.3	Partitioning Multiple cells	54
7	Performance Analysis of ORB	56
7.1	Scaling	56
7.2	Comparison to Theoretical Model	57
8	Conclusion	58

1 Introduction

The N-Body technique has been used for decades to simulate the Universe to compare theory with observations. It uses "particles" to represent objects of a certain mass and computes the forces by applying the gravity equation. As the forces operate over infinite distance it is necessary to consider all pair-wise interactions, making a naive implementation $\mathcal{O}(n^2)$. This does not scale with large particle counts.

A common solution is to partition the space, in which the particles are contained, into a set of subspaces, also called cells. These cells are then stored in a space partitioning tree data structure (SPTDS) in order to speed up the simulation with the Fast Multipole Method (FMM) to $\mathcal{O}(n)$. Building the data structure uses a significant percentage of the overall simulation time. As of now this was usually done on the CPU, not leveraging GPU acceleration. In this thesis we establish upper limits for the speedup of GPU accelerated SPTDS building algorithms over its CPU only counterpart. We then propose, implement and compare efficient methods of building the SPTDS on the CPU and the GPU.

TODO: Describe measurement results [here](#)

2 Background

Structures occurring within astrophysical simulations can be identified in real observed data[1]. Thus the simulations can help us understand the universe. Generated datasets with similar distributions to reality are used to seed the simulations. Mass field distribution are approximated with a large number of particles, where each of them represents a mass point. Numerical integration methods then compute the gravitational forces on each object. Originally developed by Joachim Stadel during his PhD at University of Washington [1], PKD-Grav is a fully parallel N-Body code which is now maintained and improved by the ICS at UHZ [2]. PkDGrav, as well as most state of the art simulation softwares, use a space partitioning tree data structure (SPTDS), which is used to partition a space into subspaces. Each cell (node), in the tree represents a subspace, and the children of the node represent the subspaces that are contained within the node's subspace. The multipole expansion is then applied to each cell of the SPTDS and the fast multipole method increases the runtime from $O(N^2)$ to $O(N)$ where N is the total number of particles. [1]

2.1 Fast Multipole Expansion

In mathematics, a multipole is a concept used to approximate the shape of a physical body. In our case the body consists of the particles contained within the subspace, or volume of a cell. The multipole expansion is the mathematical series that results from this approximation. The resulting series can simplify the original body, by choosing a suitable precision.

To keep force integration within a reasonable error margin, in most cases it is not necessary to compute all N to N interactions. Particles far away are instead considered only as a part of the volume of the corresponding cell in the SPTDS and summarized with a multipole expansion.

On a more intuitive level, one can imagine a large cluster of objects far away from planet Earth. If we wish to compute the interactions between every single object of the cluster and Earth, we will need to perform $O(N^2)$ computations. However, since the object is reasonably far away, it does not make sense to compute every interaction. Instead, we summarize the cluster into a single object and compute the gravity acting between planet Earth and the group. Consequently, objects within our own solar system would be too close to planet Earth to approximate its effects using a multipole. In such cases, it is necessary to perform all computations, however they can be hidden within a constant and do not affect the total runtime. Note that this example does not directly translate to the simulations, as a particle is not equivalent to a planet or a star. More on that in section 3.1.

2.1.1 Advantage of Cubic Bounding Volumes

From [1]: "the opening radius of a cell is given by some constant factor times r_{max} ". r_{max} is the distance between the center of mass of a cell and its most

distant corner. The opening radius correlates with the error ratio, which in turn influences how many timesteps have to be performed. Of the family of rectangular cuboids, the cube is the shape where the average distance between any point and its most distant corner point is the smallest.

2.2 ORB vs. K-d Trees

The Orthogonal Recursive Bisection method constructs a SPTDS where all resulting cells represent subspaces, or volumes, similar to cubes. K-d construction methods generate SPTDS as well, but shapes of the resulting subspaces can be strongly skewed and are thus ill suited for FMM. Both algorithms start with a single cell and iteratively construct the SPTDS by searching for axis aligned planes cutting the volume of each leaf cell of the current SPTDS into two. From the resulting two volumes, two new cells are generated and appended to the SPTDS.

In k-d tree construction, the axis lying orthogonal to the cut plane is chosen periodically, meaning all cells with an equal distance (depth) to the root of the SPTDS are cut with planes orthogonal to the same axis. In the case of ORB, the axis where a cells volume is largest is picked to compute an orthogonal plane. Therefore cells with the same depth can have cut planes of variable orientations.

2.3 Rise of the GPU

Most if not all state of the art supercomputers are hybrid, meaning each node has a CPU as well as a GPU.[3] GPU accelerated code can increase the speed of many computations, where in most cases CUDA is used to fully leverage the hardware.

The computational effort required in PKDGrav and most modern astrophysical N-Body simulations can be divided into three categories:

1. **Global Tree Building / Load Balancing:** In supercomputers, particles are distributed equally among computing nodes to leverage node level parallelization. ORB is used to generate a well suited SPTDS for FMM on a node level.
2. **Local Tree Building:** The same approach is repeated on a local level, where instead of distributing the workload among nodes, its distributed among threads. The tree from step 1. is expanded by the locally computed tree.
3. **Force Calculation and Integration:** The STPDS constructed in step 1. and 2. combined with FMM is used to compute force integration.

In PKDG each category is about one third of the total calculation time[1]. Making Global (1.) and Local Tree Building (2.) subjects to possible performance improvements, as GPU acceleration is as of now only exploited in Force Calculation and Integration(3.).

Our main objective is to improve the runtime of astrophysical simulations, or more specifically PKDGrav. Simultaneously penalties with regards to N , the total number of particles, should be kept as low as possible. Simulation accuracy benefits from large N as well as an increased number of time steps. As GPU accelerations act on a low level, hardware specific details are relevant. We want to adapt and optimize the algorithms with regards to the target computer Piz Daint.

2.4 Related Work

A number of GPU accelerated k-d tree construction algorithms exists, however the methods cannot simply be translated to ORB due to different axis choosing method. For example in the master thesis of Voglsam [4] a k-d tree is computed on the GPU and used to improve ray tracing for 3D graphics. However it makes use of a binning algorithm to find an optimal cut across cells of the same depth, which can be done in a single iteration. However in our case the axes for cells of the same depth vary.

When we subdivide ORB into its individual components as shown below, some subroutines can be computed using fairly general algorithms, which have already been implemented and optimized for the GPU in many open source libraries.

1. Cut cells on last level in tree
 - (a) Make cut plane position guess
 - (b) Count number of particles left to the cut plane
 - (c) Repeat 1. till correct plane was found
2. Partition particles
3. Repeat till desired depth was reached

Thrust by *NVIDIA* [5] has an implementation of a reduction, which can be used to perform step 1.b), a binary search is equivalent to the entire step 1. and finally it also exposes a partition interface. The library is very high level and does not allow us to control memory operations, which are very crucial to a highly efficient implementation of ORB. Furthermore CUB [6], also developed by *NVIDIA*, exposes a more low level API of the equivalent elements but the implementation is not general enough to deal with some critical performance issues. As the SPTDS grows, so do the number of leaf cells and thus the number of cells for which a cut plane has to be found in the same iteration. If we initialize a kernel for each cell, the number of kernels initialized grow exponentially, resulting in kernel invocation overheads dominating the actual computational costs.

3 Orthogonal Recursive Bisection (ORB)

Each cell in the SPTDS is enriched with the following properties:

- V_{cell} 3D volume where all corresponding particles are contained within.
- d_{cell} Number of leaf cells to be found among all subsequent cells.
- N_{cell} Number of particles encompasses in the volume of the cell.

Naturally each cell has two pointers pointing to its two children, which are used to traverse the tree, parent cell pointers important but could be added as well.

3.1 Target Data and Notation

A modified random Gaussian distributed is used to generate the target dataset, where its set to fit measured constraints of the actual universe. Thus the number of particles can set dynamically. In the simulation, a single particle has a mass of 10^8 solar masses.

- We define a particle as par_i where $i \in \{0, \dots, N\}$.
- We define the space of binary numbers with a precision of p as \mathbb{B}_p where we have $a \in \mathbb{B}_p \Leftrightarrow a \in \{0, 1\}^p$
- We define the corner coordinates of root domain with \vec{lower}, \vec{upper} where we have $\vec{lower}, \vec{upper} \in \mathbb{B}_p^3$.
- We define the coordinates of a particle par_i as \vec{x}_i for which holds $\{\vec{x} | \vec{lower} \leq \vec{x} \leq \vec{upper}, \vec{x} \in \mathbb{B}_p^3\}$. When referring to a single coordinate of a particle object we do so by writing $\vec{x}_{i,j}$. Finally the array of all particles positions in a single axis is denoted as $\vec{x}_{:,j}$.

It lies in the nature of the universe, that large clusters of particles are found in some places, whereas other areas may be vastly empty of any objects. This is a very crucial characterization and differentiates this dataset from many other application of FMM, where the distribution are more uniform.

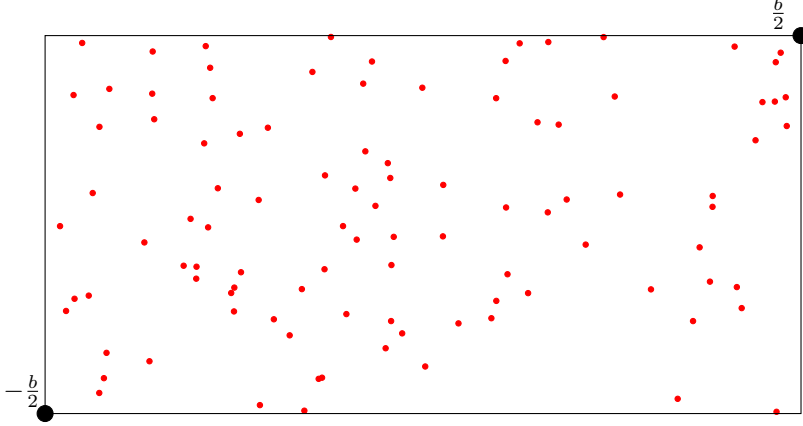


Figure 1: Uniform random distribution of 3D coordinates in cube domain projected onto a 2d plane

3.2 Memory and Workload Balancing

As mentioned we want to (1) improve the runtime of ORB whilst (2) at the same time keeping the number of particles N as large as possible. There exists some adverse effects when disregarding interactions between minimizing runtime and maximizing N .

To reduce the force integration error across all particles, some particles require more timesteps than others. This is caused by the vastly variant forces which are exerted on the particles. Whilst some follow a straight path with an almost constant velocity requiring little computational effort for great precision, others are influenced by strong gravitational poles resulting in highly curved movement paths and correlating with higher numbers of timesteps required.

We denote the workload for a particle p_i as the weighting function $w(p_i)$. The workload correlates with the number of simulations steps, that need to be computed for a single particle. In an optimally parallelized system balanced in terms of computational effort, the workload should be very similar among computing units. Where ideally the sum over all the particles considered by a thread is equal to the sum of all others. With strongly varying weights among particles, perfectly balancing the computational effort, results in a unbalanced distribution in terms of memory.

Let us consider a simple example to illustrate the point. Given the set of particles $A = \{p_1, p_2, \dots, p_{2m/3}\}$ and respectively $B = \{p_{2m/3+1}, p_2, \dots, p_m\}$. This yields a total number of particles equivalent to $N = m$. We assume that $\forall p \in A : w(p) = 1$ and $\forall p \in B : w(p) = 2$. We assign all particles from set A to process with rank 0 and the particles from B to rank 1. It follows $\sum_{p \in A} w(p) = \sum_{p \in B} w(p)$. Thus the two processors are balanced in terms of computing costs, but not in terms of memory size. In fact, process with rank 0

has $2m/3$ elements and rank 1 has $m/3$ elements. Assuming each process has a memory size of $2m/3$, then clearly this configuration is not optimal. If we were to favor memory balancing, we could assign $2m/3$ to each processor and we would be able process $(4/3) \times m$ particles in total, which is larger than the original N and therefore favorable since a dataset with a larger N could be used.

To which degree its ideal favor memory over workload balancing is difficult to answer. For now parametrize the workload using the weighting function. If we decide to completely ignore the workload balancing and solely focus on memory balancing, we can simply set $w(p_i) = 1$ for all particles.

3.3 ORB

To introduce the ORB algorithm, we consider a recursive implementation of ORB as it is easier to understand than the iterative version. We define all cell properties recursively, allowing us to describe the algorithm in detail and derive the necessary constraints.

As we initialize our recursion with the root cell alone, it follows that V_{root} is equivalent to entire volume of the input data. $d_{cell} = d$ and $N_{root} = N$ are true per definition of N and d . When a cell is cut into two child cells, we refer to the children as *leftCell* and *rightCell*. For any cell the following equations hold:

$$V_{cell} = V_{leftCell} \cup V_{rightChild} \quad (1)$$

$$V_{leftCell} \cap V_{rightChild} = \emptyset \quad (2)$$

Defining d_{cell} for each cell is especially important, since in some cases d is a non power of two number, meaning some recursion paths terminate earlier. The larger portion of leaf cells are always allocated to the left cell, as this ensures that the final constructed tree is a nearly complete binary tree.

$$d_{leftCell} = \left\lceil \frac{d_{cell}}{2} \right\rceil \quad (3)$$

$$d_{rightCell} = d_{cell} - d_{leftCell} \quad (4)$$

We can prove that the resulting tree is indeed a nearly complete binary with a case distinction:

- **Case A:** d_{cell} is divisible by two, then $d_{rightCell} = d_{leftCell}$ and the tree is complete.

- **Case B:** d_{cell} is not divisible by two, then $d_{rightCell} + 1 = d_{leftCell}$. In this case...

QUESTION: How to proof this?

Finally the number of particles encompassed in $V_{leftCell}$ and $V_{rightCell}$ are defined as follows:

$$N_{leftCell} = \min \left\{ x \in \{0, \dots, N_{cell}\} : \sum_{i=0}^x w(p_i) \geq \frac{d_{leftCell}}{d_{cell}} \times \sum_{i=0}^{N_{cell}} w(p_i) \right\} - 1 \quad (5)$$

$$N_{rightCell} = N_{cell} - N_{leftCell} \quad (6)$$

To simplify the visual and numeric explanation of the ORB algorithm, we assume $\forall i \in \{0, \dots, N\} : w(p_i) = 1$.

As the cut plane, which divides $V_{leftCell}$ and $V_{rightCell}$ is axis aligned, searching for a single value c is sufficient to find an ideal construction. Its is considered ideal if the number of particles where $x_{i,a} \leq c$ are equivalent to $N_{leftCell}$.

After a plane is found, we can divide or cut the V_{cell} in two volumes, where $V_{leftCell}$ and $V_{rightCell}$ are constrained by the original V_{cell} and the cut plane. As defined in equation 2 the volumes do not intersect, thus $V_{leftCell}$ along axis a ends at c , where $V_{rightCell}$ start at c .

3.3.1 Algorithm

Algorithm 1 is the main routine of the orb algorithm. All the volumes are described using a lower and an upper boundary point, sufficiently describing a box, thus the only input parameters required are x storing the particles positions, *lower*, *upper* and finally d . If d is equal to 1, this means the reduction must not be continued, as the target of a single subsequent leaf cell is already reached, since the cell itself is treated as a leaf cell. The stopping condition is formulated on lines 2-4. Line 6 describes a call to a method called *maxIndex()*, which essentially compares all provided values and returns the index of the maximum values. The result i provides us with the axis where the cell volume is largest. $d_{leftCell}$ is computed as described in equation 3 and its results is stored in a variable named d' . All prerequisites are met to compute the actual cut with the *cut* method which we will explain in more detail in section 3.5. The return value stored in *cut* is equivalent to the position of the cut plane along the cut axis. Using the *cut* value, the array of particles is partitioned as described in 3.6 returning the *mid* value which is the pivot index of the partitioning. Finally we can divide the original volume into two volumes as described on lines 10-13 and recursively call *ORB* method where we provide it with two slices of the x array.

Algorithm 1 The ORB main routine

```
1: procedure ORB( $x, \vec{lower}, \vec{upper}, d$ )
2:   if  $d = 1$  then
3:     return ▷ Stopping condition
4:   end if
5:    $\vec{size} = \vec{upper} - \vec{lower}$ 
6:    $i = \text{maxIndex}(\vec{size}_0, \vec{size}_1, \vec{size}_2)$  ▷ Get index of max value

7:    $d' = \lceil \frac{d}{2} \rceil$ 
8:    $cut = \text{cut}(x, \vec{lower}_i, \vec{upper}_i, i, \frac{d'}{d})$  ▷ Find cut plane
9:    $mid = \text{partition}(x, split, axis)$  ▷ Partition particles

10:   $\vec{upper}' = \vec{upper}$ 
11:   $\vec{upperChild}_i = cut$ 
12:   $\vec{lower}' = \vec{lower}$ 
13:   $\vec{lowerChild}_i = cut$ 

14:  ORB( $x_{0:mid}, \vec{lower}, \vec{upper}', d'$ )
15:  ORB( $x_{mid:len(x)}, \vec{lower}', \vec{upper}, d - d'$ )
16: end procedure
```

3.4 By Example

A sample dataset is proposed to help with visual and numerical explanations of the algorithm. All particles in the sample dataset are assumed to have a weight of 1.

We explore the algorithm visually using the provided example dataset:



Figure 2: Example distribution with $N = 7$



Figure 3: Example particles with ORB at recursion depth 0



Figure 4: Tree with ORB at recursion depth 0

Depicted in figure 3 & 4 is the SPTDS after initialization: there is only one

$cell_1$ which is root and V_{cell_1} encompasses the entire domain which in this case is the rectangle ranging from 0 to 1 in both the x and y axis. d_{cell_1} is equivalent to 3 and N_{cell_1} is 7.

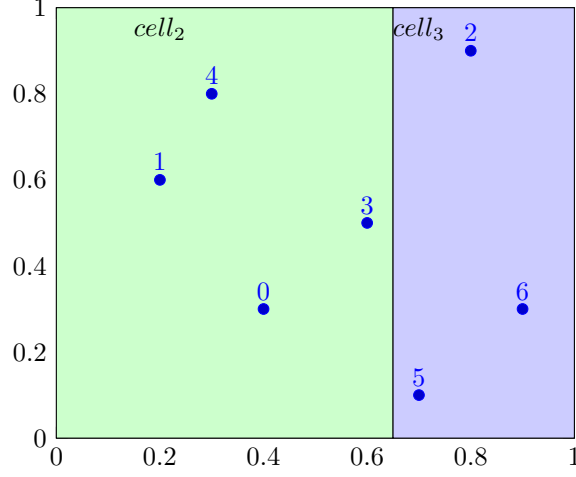


Figure 5: Example particles with ORB at recursion depth 1

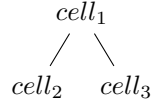


Figure 6: Tree with ORB at recursion depth 1

A cut plane, which translates to a cut line in simplified 2D example is constructed and $cell_2$ and $cell_3$ are generated accordingly.

5 & 6 Depict a SPTDS of depth 2, or two levels, where a cut plane was found dividing V_{cell_1} into V_{cell_2} and V_{cell_3} , $d_{cell_2} = 2$ and $d_{cell_3} = 1$ subsequently are calculated, meaning $cell_3$ is considered a leaf cell. Thus we can compute N_{cell_2} using a simplified version of equation 5 since we assume all weights are equal 1: $N_{cell_2} = \lfloor \frac{d_{cell_2}}{d_{cell_3}} * N_{cell_1} \rfloor = \lfloor \frac{2}{3} * 7 \rfloor = 4$ subsequently $N_{cell_3} = 3$.

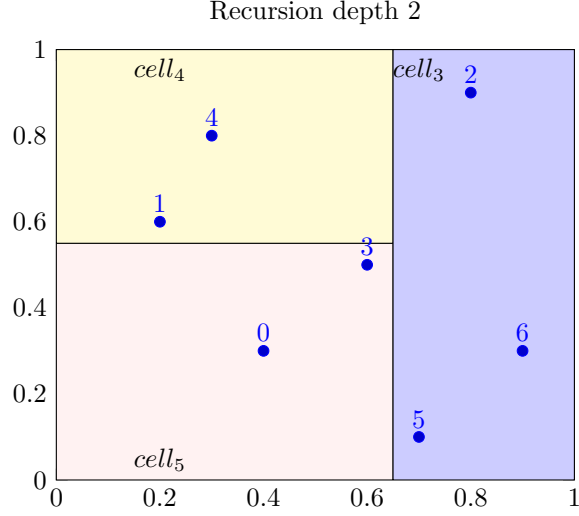


Figure 7: Example particles with ORB at recursion depth 2

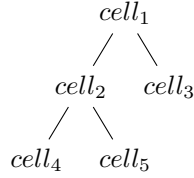


Figure 8: Tree with ORB at recursion depth 2

Again the procedure is repeated for $cell_2$ but not for $cell_3$ as there recursion is terminated by the stopping condition. Finally we end up with $cell_4$ and $cell_5$ where the stopping condition is met as well. The resulting SPTDS has 3 leaf cells and partitions the space into three subspaces.

3.5 Root finding

The *cut* algorithm takes an array of particle positions x , an *axis*, *left* and *right* boundaries and a *percentage*. Its goal is to return a position along the cut axis such that the particles less or equal to the a *cut* value are equivalent to the percentage multiplied by the length of the x array.

The problem is related to a selection algorithm and a quick select[7] could be used. However quick select has a worst case runtime of $O(n^2)$. Furthermore median algorithms could be explored as well, but they are mostly approximation algorithms with bad worst case runtimes. We can reformulate the problem as a root finding problem. To do so we define a function $f(c)$ which evaluates

the number of elements smaller than the cut value minus half the number of total particles. When the function evaluates to zero, we have found an axis aligned cut plane, where exactly half the particles are located on the left side. There exist many different solvers for the root-finding problem, but the most stable and easiest to implement is the bisection method. Some solver combine approximate solver with the stable bisection method to generate fast but stable root finding methods. Such algorithms could be explored in later work.

Since we are operating on binary numbers with a limited precision, the bisection method is guaranteed to terminate after p steps.

3.5.1 Bisection Method

Initially an estimation for a cut is made, in this case the exact middle of the domain boundaries. Some median finding algorithms use improved guessing to speed up the process, again a method which could be leveraged in later work.

Because the maximal number of iterations is known, a loop can be used as seen in algorithm 2 on line 4. A cut is then computed (line 5), which in the first iteration is the center position between left and right. We then check weather the stopping condition has already been reached (line 7-9), meaning the cut lies within α points of the ideal result. If it was not reached the boundaries can be improved as follows: In case there are too many elements left of the cut, we know that the cut position was chosen too far to the right and we can be sure that the ideal cut must be left of the current guess. Therefore we adjust the boundaries, in this case we set *right* equal to *cut* as seen on line 13. The analog concept can be applied in the other case (line 11).

Algorithm 2 Bisection Method

```
1: procedure CUT( $x, left, right, axis, percentage$ )
2:    $nLeft = 0$ 
3:    $cut = 0$ 

4:   for  $k \in 0, \dots, p$  do
5:      $cut = (right + left)/2$ 
6:      $nLeft \leftarrow \text{sum}(x_{:,axis} < cut)$  ▷ Counting particles left of cut

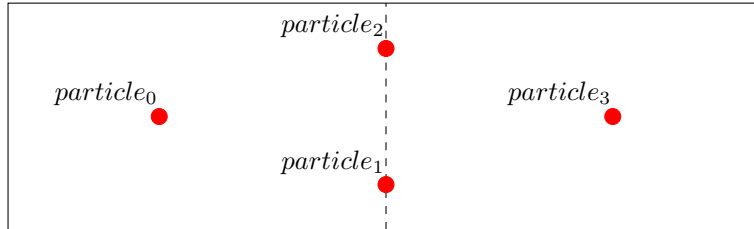
7:     if  $\text{abs}(nLeft - \text{len}(x)) * percentage < \alpha$  then
8:       Break ▷ Stopping condition
9:     end if

10:    if  $nLeft \leq \text{len}(x) * percentage$  then
11:       $left = cut$ 
12:    else
13:       $right = cut$ 
14:    end if

15:  end for
16:  return  $cut$ 
17: end procedure
```

3.5.2 Edge Cases

Let us consider the following example particle distribution where $particle_1$ and $particle_2$ have identical x coordinates.



In this case there exists no ideal cut in the x (horizontal) axis. As either there is 1 particle to the left or 3, but no cut can result in 2 particles to the left. If we keep adding particles along the same line, the method performs even worse. However since the algorithm uses a deterministic for loop, the algorithm will terminate anyways.

QUESTION: Write this? seems too vague

3.5.3 Runtime Analysis

On line 5, the number of particles to the left of the cut plane are summed. As the array is unordered, all coordinates of one axis need to be read once, resulting in a runtime of $O(N)$. All other operations inside the for loop can be computed in a negligible constant time. The loop itself is repeated p number of times, where p corresponds to the precision of x coordinates.

To proof that the loop concludes after p iterations we assume integer numbers. With each iteration the range of possible solutions is divided by two. The same goes for integer numbers, each time a bit is removed, the size of the range of numbers which can be represented is reduced by two. Thus after p iterations the precision limit is reached and the cut cannot be improved.

3.6 Partition Algorithm

We want to continuously update the array storing the positions of the particles and in the words used before, have a direct correlation between x_i and i . This enables grouping particles within a cell in a fixed range of two indexes of the particles array. The advantages of this are two fold: Access all particles within a cell in constant time, manipulate particles within a cell using a slice of the array.

As seen in the pseudo code in algorithm 3, the algorithm looks for a pair of particles, where for both the coordinate along the relevant axis are on the wrong side of the cut plane. In this case the particles can be swapped (line 8) resulting in the correct position of both particles. We refer to [7] for a correctness proof and a more detailed explanation of the algorithm.

Algorithm 3 Partition Method

```

1: procedure PARTITION( $x, cut, axis$ )
2:    $i = 0$ 
3:   for  $k \in 0, ..N - 1$  do
4:     if  $\vec{x}_{k,axis} \leq cut$  then
5:       while  $\vec{x}_{i,axis} \leq cut$  and  $i < N$  do
6:          $i = i + 1$ 
7:       end while
8:        $x_i, x_k = x_k, x_i$ 
9:     end if
10:  end for
11:   $x_i, x_{N_j-1} = x_{N_j-1}, x_i$ 
12: end procedure

```

A runtime of $O(N)$ can be derived, as the algorithm iterates over all particles once. Since we need to touch each element at least once to partition the entire array there exists no better method.

Lets apply the algorithm to our running example. We start with our initial array of particles as follows:

x	y	id
0.4	0.3	0
0.2	0.6	1
0.8	0.9	2
0.6	0.5	3
0.3	0.8	4
0.7	0.1	5
0.9	0.3	6

We then partition the particles with a cut in the x axis set to 0.65 as seen in figure 5.

x	y	id
0.4	0.3	0
0.2	0.6	1
0.3	0.8	4
0.6	0.5	3
0.8	0.9	2
0.7	0.1	5
0.9	0.3	6

Finally we partition $cell_2$ into $cell_4$ and $cell_5$ with the cut position 0.55 along the y axis as seen in figure 5. We end up with the following array:

x	y	id
0.4	0.3	0
0.6	0.5	3
0.3	0.8	4
0.2	0.6	1
0.8	0.9	2
0.7	0.1	5
0.9	0.3	6

Note how all particles from $cell_4$ are contained in the range of 0-1. The particles of $cell_5$ in 2-3 and finally the particles contained in the volume of $cell_3$ can be found in 4-7.

4 Theoretical Analysis

The main goal of this thesis, is to improve the runtime of the ORB algorithm by leveraging the graphics processing unit over the central processing unit. Before implementing, it makes sense to explore if and how strong the performance benefits could in theory manifest itself. For this purpose we develop a simplified model, which we use to compare the CPU version against two different GPU variants. As GPU programming is very low level, the knowledge gained while developing the model also builds a solid theoretical foundation.

4.1 General Memory Model

For a consistent terminology of a computer, we briefly establish a general computing model, which can be applied to most modern high performance systems. We are looking at a single node and its hardware components, a supercomputer may have thousands of these computers linked together where different bandwidth are seen between individual nodes.

Both the CPU and the GPU have their own memory which are connected by a data link, where the connections are bound by the memory bandwidths. We name the capacity of the CPU memory bandwidth B_{CPU} and the GPU memory bandwidth B_{GPU} . There is a separate data link between the CPU and the GPU memory, which is commonly referred to as PCI express or NVLink (for modern *NVIDIA* GPU's, in our model we use the term I_{GC} . In systems with multiple CPU there is a link between the individual CPU's which we denote as I_{CC} . Finally we call the link between GPUs as I_{GG} .

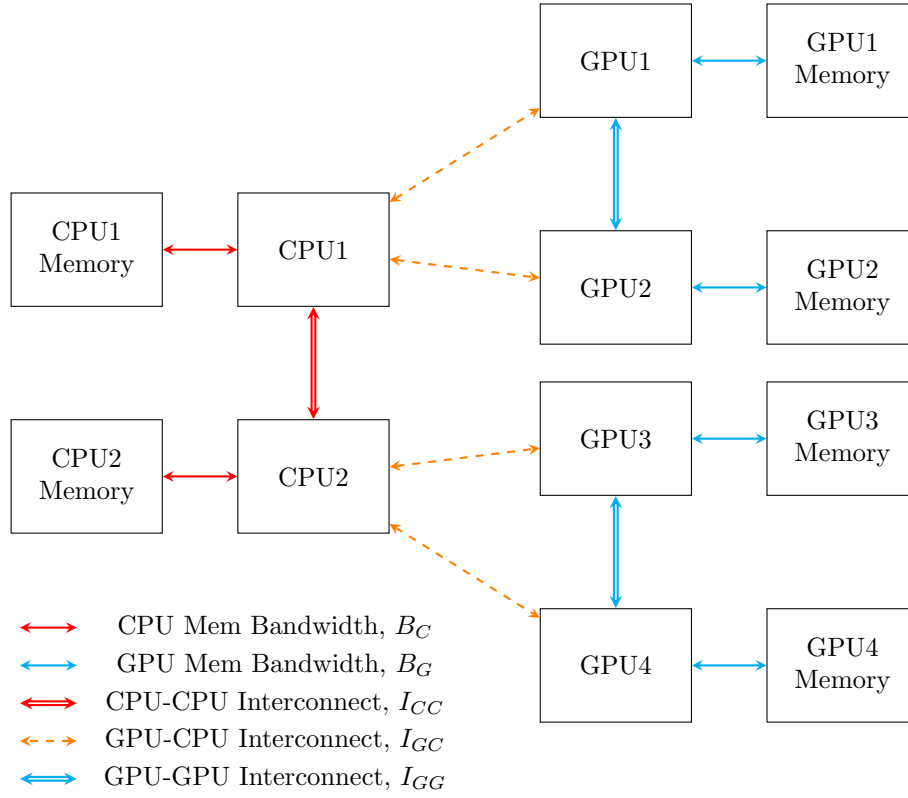


Figure 9: Illustration of the universal computing model

4.2 Supercomputers

For Piz Daint, Summit and Eiger we have collected all relevant hardware metrics and compiled it into the table seen in table 1. Piz Daint and Eiger were chosen, because we have the possibility to test the Code on its systems. Furthermore we include Summit as its is at the time of writing this thesis, one of the most capable supercomputers in the world.

Constant	Piz Daint [8]	Summit[9]	Alps (Eiger)
# Nodes	5704	4608	1024
# CPUs	1	2	2
CPU Model	Intel E5-2690 v3 [10]	IBM POWER9	AMD EPYC 7742[11]
CPU Mem.	64 GB	256 GB	??
B_C	68 GB/s	170 GB/s	204.8 GB/s x 2
I_{CC}	-	64 GB/s	??
Base GHZ_C	2.9 GHZ	4 GHZ	2.25 GHZ
Max GHZ_C	3.8 GHZ	4 GHZ	3.4 GHZ
# Cores	12	22	64
Architecture	Haswell	POWER9	AMD Infinity Architecture
SIMD	AVX2	?	AVX2
# GPUs	1	6	0
GPU Model	NVIDIA P100 [12]	NVIDIA V100s [?]	-
GPU Mem. Cap.	16 GB	16 GB x 6	-
B_G	732 GB/s	900 GB/s x 6	-
I_{GC}	32 GB/s	50 GB/s x 6	-
I_{GG}	-	50 GB/s	-
GPU Tflops	9,3	16.5	-
# CUDA Cores	3584	5120	-

Table 1: Datapoints of Supercomputers

4.3 Roofline Performance Model

In a first step we determine whether our computations are bound by memory bandwidth or the actual performance of the computing unit. The most costly computation is line 6 of the cut method (figure 2) with a runtime of $O(32 \times N)$. Oftentimes when an algorithm iterates over a large dataset performing only very little calculations on its individual elements, the limiting factor is the memory. To support this claim we make roofline models for all three systems and check whether they are really memory bound. A roofline model compares arithmetic intensity in flops per byte against the actual performance of the computing chip in flops. Therefore we need to gather all the relevant data first.

4.3.1 Estimating Flops

For modern hardware, it is fairly uncommon to release flops (floating point operations per second) values. Steadily evolving SIMD instruction sets result in varying performance for different implementation details, which in turn are compiled into different assembly instructions. Depending on the algorithm, implementation details and compilation flags the c++ compiler tries to compile ideal assembly instruction sets. SIMD instructions can only be used when there is a contiguous memory access, thus in some cases a poor memory layout choice may lead to a much lower flops.

For most modern CPU chip architectures AVX is the fastest SIMD instruction set available. The common AVX2 enables the processing of 8 floating point operations per instruction with a CPI (cycles per instructions) of 0.5. The CPI can also vary depending on chip architecture, but to our knowledge all relevant CPU's from figure 1 indeed support a CPI of 0.5 along with AVX2. A lower CPI results in a higher efficiency, as several instructions can be completed in a single cycle.

We define in equation 4.3.1, a function to estimate the number of gigaflops for a given hardware. We define GHZ as the gigahertz which can be reached by the CPU, meanwhile NF is the number of floats which can be processed simultaneously using AVX. CPI is as mentioned the cycles per instructions for the AVX instruction set. Finally we have np which is the number of processors, where we assume a perfect parallelization, meaning 100% of the code can be parallelized.

$$gflops = GHZ * NF * CPI^{-1} * np \quad (7)$$

Since we have peak flops benchmarks available for the *NVIDIA* P100 and V100s we do not need to make any estimations on the GPU side.

4.3.2 Estimating Arithmetic Intensity

Arithmetic intensity is measured in FLOPS per byte or the number of floating point operations which are computed per byte loaded from memory. The count left part from the cut algorithm (line 6) in figure 2 can be translated to the following isolated c++ code:

Listing 1: Counting the particles left of a cut plane

```
1 for(auto p= startPtr; p<endPtr; ++p) nLeft += *p < cut;
```

Where p is a C-style array which stores the particles position and $nLeft$ stores the number of particles which are smaller than cut . We can list the operations per loaded float:

1. Compare particle to cut
2. Add result to $nLeft$
3. Increment pointer p
4. Compare pointer p with $endPtr$

Which results in a total of 4 operations. Since a single float is stored using 4 bytes, this computes to 1.0 operations per byte or an arithmetic intensity of four.

Estimating the arithmetic intensity for the GPU is a lot more complicated as it can vary a lot depending on the specific implementation details. But for now we will just assume the same arithmetic intensity as we have had for the CPU.

Note that SIMD instructions do not influence the arithmetic intensity, as the number of floating point operations per byte remains the same. We simply perform a set of operations concurrently. What is however influenced by AVX, is the maximum number of GFLOPS which can be processed by the hardware. If we were to ignore AVX, the algorithm would clearly not be bound by memory, as its performance would lack behind the memory bandwidth. For this reason it is important to consider SIMD instructions.

4.3.3 The Plots

Lets us compute the maximally achievable gflops for Piz Daint. To do so, we plug in the corresponding datapoints from figure 1 into the formula 4.3.1 for $np = 1, 2$ and 3 .

$$3.8 * 8 * \frac{1}{0.5} * 2 = 121.6gflops \quad (8)$$

$$3.8 * 8 * \frac{1}{0.5} * 4 = 243.2gflops \quad (9)$$

$$3.8 * 8 * \frac{1}{0.5} * 8 = 486.4gflops \quad (10)$$

The results are depicted as horizontal lines in 10. The memory bandwidth is plotted as a line with an equivalent slope. Finally we add a dotted line representing the arithmetic intensity of the Count Left procedure. To interpret the roofline mode, one has to follow the dotted line starting from the bottom. Weather it intersect with line representing the maximal performance or the memory bandwidth first, gives an indication weather the program is performance or memory bound.

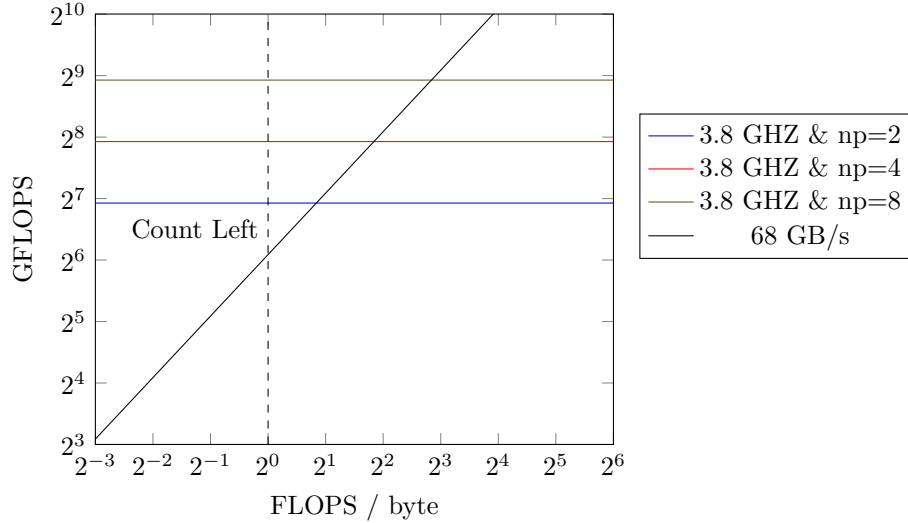


Figure 10: Roofline Model for Piz Daint CPU

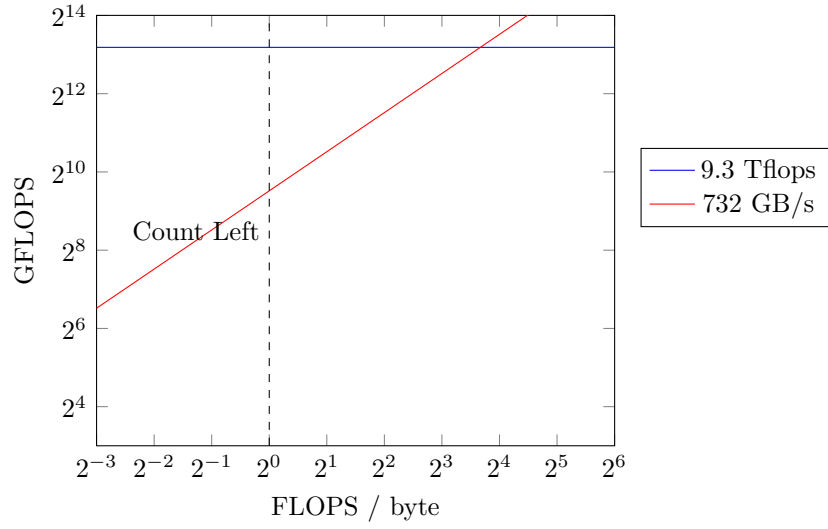


Figure 11: Roofline Model for Piz Daint GPU

As it can be seen in figure 10 and 11 both the GPU and CPU version running on Piz Daint are memory bound, but due to its much higher memory bandwidth limit, a GPU version should be favored.

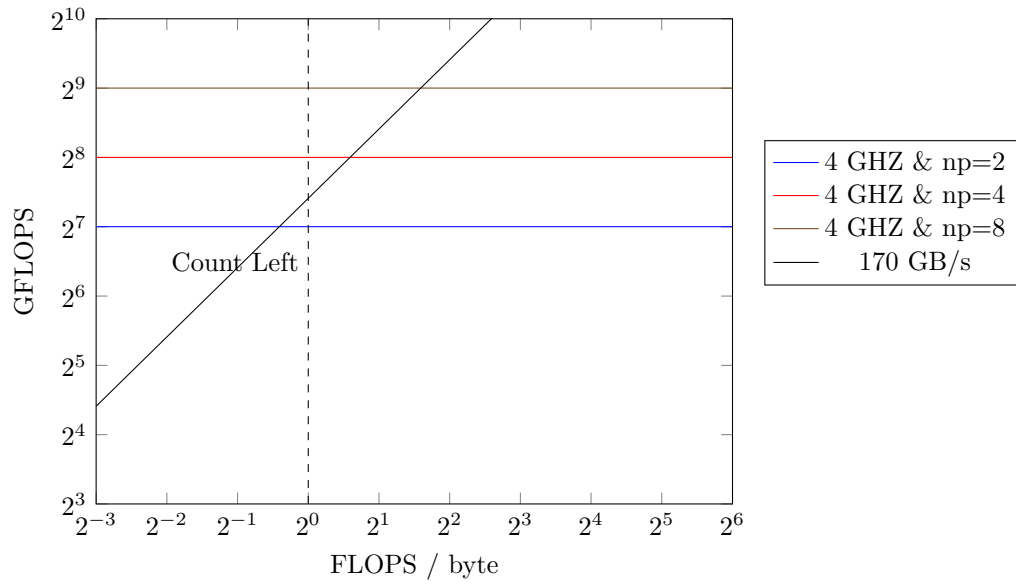


Figure 12: Roofline Model for Summit

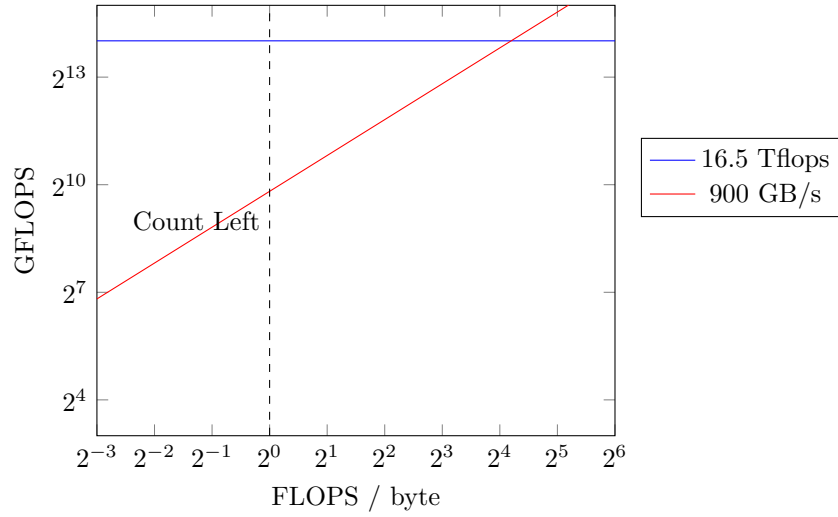


Figure 13: Roofline Model for Summit GPU

As visible in figure 12 and 13 the same applies to Summit. Its notable how the memory bandwidth only becomes the limiting factor when assuming a parallelization with 4 processors.

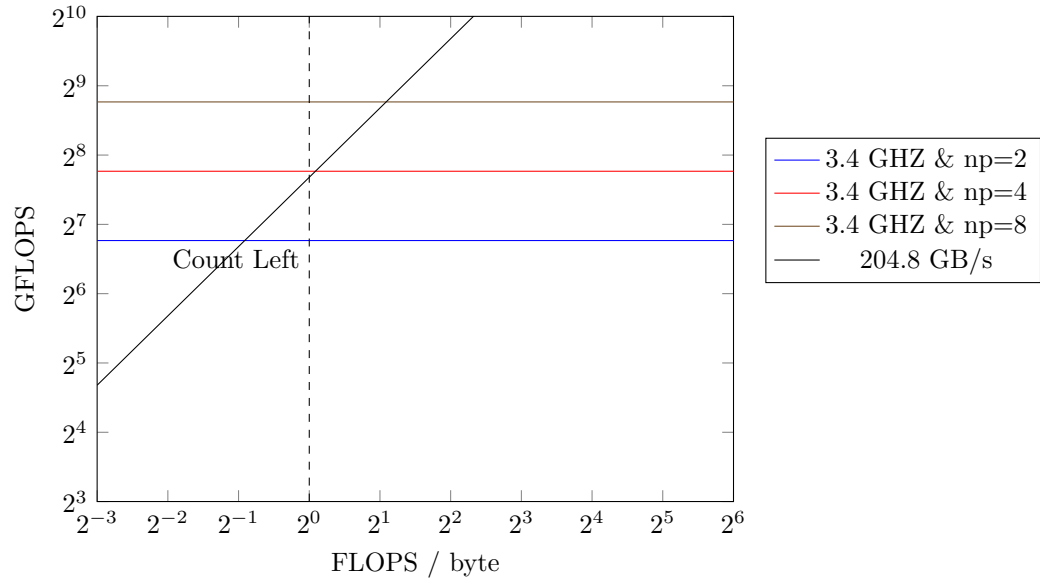


Figure 14: Roofline Model for Alps

Finally alps, which has the highest memory bandwidth compared to its flops, also becomes limited by the bandwidth after using more than four processors.

4.3.4 Empirical Verification

We consider a minimal C++ code to verify the theoretical model. We use the -S flag along with the g++ compiler to generate assembly code from the c++ source code and make sure that AVX commands are enabled. For testing purposes all entries in the array are set to random values between 0 and 1, and cut is set to 0.5.

In a first test, we do not enable AVX, but turn on O3. The generated assembly code is depicted in listing 2.

Listing 2: Reduction Assembler Code without AVX

```

1 .L18:
2     movups    (%rax), %xmm0
3     addq     $16, %rax
4     cmpltps   %xmm2, %xmm0
5     psubd    %xmm0, %xmm1
6     cmpq     %rdx, %rax
7     jne      .L18

```

psubd is a packaged instruction, meaning it already uses some form of SIMD instructions. The command is used in the MME and later the SSE2 instruction set. Since we can observe that the *xmm* registers are used, we know its a SSE2 instruction.

If we additionally set the compile flag -march=native we end up with the instructions as depicted in listing 3.

Listing 3: Reduction Assembler Code with AVX2

```

1 .L19:
2     vmovups   (%rax), %ymm3
3     addq     $32, %rax
4     vcmpltps  %ymm2, %ymm3, %ymm0
5     vpsubd    %ymm0, %ymm1, %ymm1
6     cmpq     %rdx, %rax
7     jne      .L19

```

We can identify *vmovups*, *vcmpltps* and *vpsubd* which are AVX commands. Since they are using the *ymm* instead of the *xmm* registers, we know these are AVX2 and not AVX commands.

QUESTION: [Look at this with doug](#) The following results were obtained on 2^{27} with a thread count of 16 on Piz Daint. The reduction takes 72010 microseconds. This means we have a throughput of $2^{27}/10^9 * 10^6/72010 = 1.8gflops$. This of course does not lie anywhere near the theoretical maximum, which even for a single processor ($np = 1$) is $3.8 * 8 * 2 * 1 * 16 = 972.8gflops$.. This is a strong indication, that already with a single processor solutions, we are in the realm of bandwidth limited algorithms. This effect of course, would become even stronger when including parallelism.

4.4 Runtime Estimates

We construct runtime estimation for the CPU version and two different GPU implementations. In the first GPU implementation *GPU Counting*, we simply send an array

of particles to the GPU, where counting is then performed. In a more advanced variant *GPU Counting and Partitioning* we send the particles from the CPU to GPU only once and perform the partitioning on the GPU as well.

We assume a precision p of 32, which is general standard for both integers and floats and its a sensible assumption for astrophysical simulations. The total storage required for all particle positions is $32 \times 3 \times Nbits = 4 \times 3 \times NBytes = 12 \times NBytes$. Furthermore we assume $d = 1024$ and $N = 10^9$.

4.4.1 CPU Version

Whilst the number of cells increases each time a level is added, simultaneously the number of particles to be iterated over . Thus the amount of cells which are to be iterated cancels out and we end up with the size of all particles s divide by the memory bandwidth B_C .

Each time the leaf cells of the SPTDS are split into two child cells and appended to the tree, we end up with twice the number of leaf cells. Considering we want to end up with d leaf cells in the end, we need to perform the cut algorithm $\lceil \log_2(d) \rceil$ times for all leaf cells.

Consequently we sum over all $\lceil \log_2(d) \rceil$ iterations, where in each iteration a cut has to be found for i cells. The cut itself is found in p iterations, where the actual array of particles is each time we grow the SPTDS by a level. As only the memory bandwidth is considered, we divide the size of the particles to be iterated over.

$$\sum_{i=1}^{\lceil \log_2 d \rceil} i \times p \times \frac{s}{B_C} \quad (11)$$

i cancels out from the equation and we end up with:

$$\sum_{i=1}^{\lceil \log_2 d \rceil} p \times \frac{s}{B_C} \quad (12)$$

Finally we can simplify to:

$$\lceil \log_2(d) \rceil \times \left(p \times \frac{s}{B_C} \right) \quad (13)$$

4.4.2 GPU Counting

The Equation 14 for the GPU is similar, the only difference being that we use the GPU memory bandwidth B_{GPU} instead of the CPU bandwidth. Furthermore we have to consider the time it takes to send the data from the CPU to the GPU. This adds the terms size divided by CPU to GPU memory bandwidth denoted as I_{GC} . Finally we also need to load the data from the CPU memory to the CPU before we are able to send it.

$$\lceil \log(d) \rceil \times \left(p \times \frac{s}{B_G} + \frac{s}{I_{GC}} + \frac{s}{B_C} \right) \quad (14)$$

4.4.3 GPU Counting and Partitioning

We can also consider performing the partitioning on the GPU, this means that there is no need to send data from the CPU to GPU each time we want to find a cut, allowing us to reduce the costly overheads. This means that we are able to move the corresponding terms out of the brackets.

$$\lceil \log(d) \rceil \times \left(p \times \frac{s}{B_G} \right) + \frac{s}{I_{GC}} + \frac{s}{B_C} = t \quad (15)$$

4.4.4 Plugin Values

Let us plugin the values from figure 1 into the corresponding formulas 13, 14 and 15 for Piz Daint. The naive implementation yields the following speeds for the naive CPU implementation:

$$\lceil \log(1024) \rceil \times \left(32 \times \frac{12GB}{68GB/s} \right) = 56.47841s \quad (16)$$

And the corresponding GPU implementation:

$$\lceil \log(1024) \rceil \times \left(32 \times \frac{12GB}{732GB/s} + \frac{12GB}{32GB/s} + \frac{12GB}{68GB/s} \right) = 10.762s \quad (17)$$

This yields in a speed-up of:

$$\frac{56.47841}{10.762} = 5.24794 \times \quad (18)$$

When considering the GPU Counting and Partitioning we end up with: And the corresponding GPU implementation:

$$\lceil \log(1024) \rceil \times \left(32 \times \frac{12GB}{732GB/s} \right) + \frac{12GB}{32GB/s} + \frac{12GB}{68GB/s} = 5.79802s \quad (19)$$

This yields in a speed-up of:

$$\frac{56.47841}{5.79802} = 9.74097 \times \quad (20)$$

4.5 Conclusion

The computations for Eiger and Summit can analogously to the ones from Piz Daint. For Eiger we cannot make any assumption for a GPU version as its CPU only system. All the results are plotted in 15. As expected the GPU Counting outperforms the CPU version on both hybrid (CPU and GPU available) super computers. Furthermore GPU Counting and Partitioning yields more performance improvements, but the speed up is less drastic. As the model may differ from reality depending in implementation and compilation details, these results cannot be taken for granted and we will compare the results with actual empirical observations in section 7.

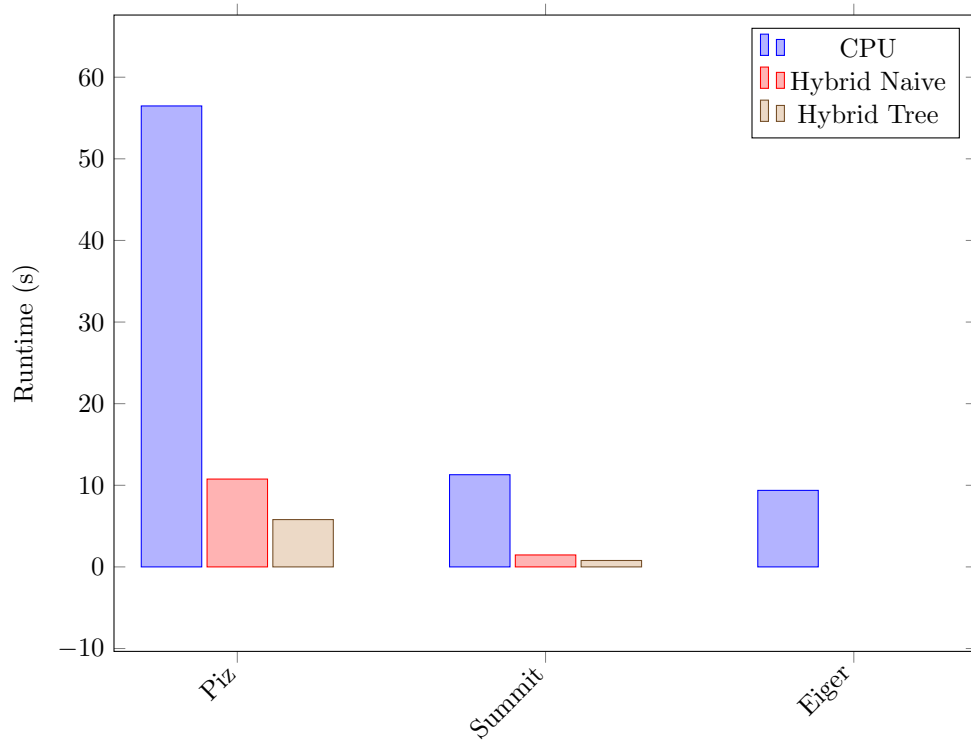


Figure 15: Execution times of different strategies

5 CPU Implementation

In this section we describe the stand-alone implementation of the ORB algorithm using CPU acceleration only. Building on the CPU version, we later introduce GPU accelerated alternatives to selected parts of the code.

We develop somewhat separately from the PKDGrav codebase in order to keep the project smaller and more easy to understand. However we make use of the machine dependent layer (MDL) from PKDGrav, which is used to distribute the workload among processors. This reduces the complexity of the implementation and we do not have to worry about communication methods and CPU parallelism too much. Furthermore it makes integration into PKDGrav simple. Note that processors might be located in different nodes, but MDL is general enough to abstract node and processor level communication and parallelization away.

We refrain from describing each part of the code in detail, however any details crucial to performance are described. Furthermore, we will pay special attention to the CUDA kernels and also memory management strategies for several reasons: minor implementation details can have a big impact upon the performance and increased replicability.

For this project we will use C++ 17 along with CUDA 11.3, OpenMPI 4.1.0 and Blitz++ 1.0.2 [13].

5.1 MDL and PKDGrav

MDL, the machine dependent layer provides an abstraction layer around the parallel primitives introduced by PKDGrav. The Master layer is responsible to coordinate the flow of the program. It does so by calling the PST (Processor Set Tree), which distributes the tasks and gets processors to work on them. As the name suggests, the PST is organized as a binary tree where intermediate nodes contains a pointer to a consecutive processor as well as a next lower node. Each leaf cell of the graph then correspond to an individual processor.

Parallel processes are dispatched by descending the PST until they reach all processors. There computations are performed and the results are combined by passing them back up the PST until they reach the master. A parallel process can be programmed by implementing the service interface as explained in section 5.5.

Each processor has its own local data, which can be accessed by calling *pst->lcl*.

5.2 Particles

The particles array has ownership over all particle objects. It has a (N,3) shape where each row represents a single particle. To enable coalesced memory access when iterating over a single axis i.e. $x_{:,axis}$, we choose a row major storage order, over a column major storage order. Even though column accesses patterns i.e. $x_{i,:}$ are used as well, they happen less frequently. Most importantly there are cases where entire section of a row have to copied, which is by magnitude faster with row major storage.

We use C-style arrays interchangeably with Blitz++ arrays to store all particles. Blitz++ is an open source wrapper class around C-style arrays, which helps with pointer management and can speed up the debugging process by keeping track of the array boundaries. Furthermore it features array slicing and does its own memory management. In certain cases we need to access the actual C-style array, which can be accessed with *array.data()*. For example when iterating over the Blitz++ array,

we have observed compiled assembly code which does not reflect the latest SIMD instructions.

The particle data is loaded into the local data storage of each processor, meaning each processor owns a different unique set of particles.

5.3 Cell

The cell class is a structure keeping track of the fundamental cell information. In essence it is the analogue to the concept of the *cell* which we have already introduced in section 3.

Quickly building and accessing elements of the SPTDS generated by ORB requires a suitable data structure. Instead of a tree with pointers, we can use a heap since the necessary conditions for a nearly complete binary tree are met as seen in section ?? . Elements within a heap can be accessed and added in constant runtime. Depicted in figure 16 is the finished SPTDS of the sample dataset stored as a heap.

The SPTDS and thus the cells are constructed on the master alone, subsets cells are only distributed along the SPTDS when launching certain services.

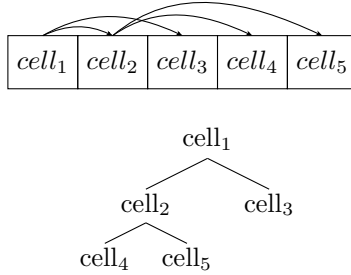


Figure 16: Tree as heap

An array can be used as an efficient storage medium for the heap, since in our case d can be set at compile time and the maximum number of cells in the SPTDS can be derived from that. Meaning the array can be allocated statically. Because MDL communicates data as arrays between the threads, another advantage of the array storage is the absence of any costly data structure conversion.

Since the STPDS is only a nearly complete binary tree and not a complete binary tree, we have to be careful when iterating over the levels of three, as the very last level is not filled entirely with cells. But because d is known, many attributes of the SPTDS can be determined deterministically.

5.3.1 Heap Conditions

As with all heap data-structures we have the following conditions:

- The root cell has index 1
- The left child of cell at index i can be accessed at the index $i \times 2$.
- The right child of cell at index i can be accessed at the index $i \times 2 + 1$.

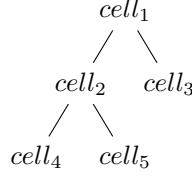


Figure 17: Tree with $d = 3$

We can derive the following constraints given the number of leaf cells d .

- We can compute the depth of the tree with: $\lceil \log_2(d) \rceil$
- The number of leaf cells on the second last level is given by $2^{\text{depth}} - d$
- There are exactly $2 \times d - 2^{\text{depth}}$ items (which must all be leaf cells) on the last level.
- The total number of cells are $2 \times d - 1$.

In the example tree depicted in figure 17 we indeed observe that the depth corresponds to: $\lceil \log_2(d) \rceil = 2$, the number of leaves on the second last level is $2^2 - 3 = 1$ and the number of leaves on the last level are $2 \times 3 - 2^2 = 2$. Finally the total number of cells is $2 \times 3 - 1 = 5$.

5.3.2 Class

Among the fields in the cell are the following variables:

- *id*: Unique identification of the cell instance. Corresponds to its index in the heap array plus one. The plus one comes from the different indexing for heap constraints, where we start with 1, and the classic array indexing starting with 0. It can be used to compute indices of both child cells and parent cells using the above shown formulas.
- *nLeafCells*: Equivalent to d_{cell} and depicts the number of leaf cells to be found in all of its successors. Its used when building the tree, to track whether a cell needs to be split or not.
- *lower*: 3D point coordinate representing the lower boundary corner of the 3D volume V_{cell} which is encompassed by this cells.
- *upper*: Represents the upper boundary corner of the volume.
- *cutAxis*: Axis where the cell is to be cut.

To support and simplify common operations, a further abstraction layer is given with the following functions:

- *setCutAxis()*: Evaluates the proper *cutAxis* given the *upper* and *lower* variables.
- *cut()*: Cuts the cell and returns two new cells.
- *lower*: 3D point coordinate representing the lower boundary corner of the 3D volume V_{cell} which is encompassed by this cells.
- *getLeftChildId()*: Returns the ID of the left child cell, meaning its volume is in the cut axis smaller than the cut plane.
- *getRightChildId()*: The equivalent for the right child.

5.4 Mapping Cells to Particles

As SPTDS heap is stored on the master but the particles are distributed among the local storages of the processors, its is crucial to know which particles are encompassed in the volume of which cell. Thus we introduce a data structure to keep track of the relation between cells and its particles. Thanks to the partitioning algorithm all particles belonging to a single cell can be found in a consecutive section of the particles array. Thus we can use a 2D array, which again is stored in the local data of the processor, where for each cell the corresponding range is stored as a tuple. Making the size of the array equivalent to the total number of cells $2 \times d - 1$ multiplied by two. Naturally the data structure needs to be updated across all processors whenever we a partition is performed.

5.5 Services

A service implementation consists of a *Service()* function as well as a *Combine* function. The service is handed two a *PST* class which stores information about the process location as well as the local data. Furthermore two void pointers storing providing storage for the input and output data called *in* and *out* along with their respective sizes. The reason void pointers are used, is the actual data structure of the input is variable for each service and this solution provides enough flexibility. Data stored inside the void pointers can be converted into their respective classes using casting. The finished results of the computations are then stored inside the output pointers.

In the *Combine()* function two output void pointers *vout* and *vout2* and the sizes of the underlying arrays are given as parameters. A combination strategy of the elements can be chosen and implemented, where the results is to be stored in *vout*. For example if we were to implement a service which returns the sum of all particles, the combine function would simply store the sum of *vout* and *vout2* in *vout*.

5.5.1 Init and Finalize Service

All data stored as local data, which is not managed by Blitz++, must be allocated when initializing and after the program has finished, the memory is freed to avoid any leaks. Allocation and freeing is not free of cost, thus it makes sense to reuse memory whenever possible.

The initialize service is also responsible to load the particle data into the local memory of each PST.

5.5.2 Count Left Service

As mentioned counting the elements smaller than the cut position, or as we call it the count left procedure is computationally intensive. For this reason we introduce a service which along with PST distributes the task among the processors. For the *Service()* function which performs the actual computations, *in* points to an array of cells and as we know the length of the array, we can simply iterate over all of its elements and cast them to the correct cell data structure as shown in listing 4 on line 1-2. In case the *foundCut* flag of cell was set to true, we know that the ideal cut was already found, thus we can continue with the next cell (line 4-6). We then read the begin and end indices of the particles array corresponding to the range of objects contained within the cells volume and respectively make a slice of the local particles array (lines 7-11). We then start and end pointers and obtain the cut position to be

tested (lines 13-17). We then count the number of particles smaller than the cut and write the result in the proper index of the output arrays (lines 18-23).

Listing 4: Part of the Count Left Service() method

```

1 for (int cellPtrOffset=0; cellPtrOffset<nCells; ++cellPtrOffset){
2     auto cell = static_cast<Cell>*(in + cellPtrOffset);
3
4     if (cell.foundCut) {
5         continue;
6     }
7     int beginInd = pst->lcl->cellToRangeMap(cell.id, 0);
8     int endInd = pst->lcl->cellToRangeMap(cell.id, 1);
9
10    blitz::Array<float,1> particles =
11    pst->lcl->particlesAxis(blitz::Range(beginInd, endInd));
12
13    float * startPtr = particles.data();
14    float * endPtr = startPtr + (endInd - beginInd);
15
16    int nLeft = 0;
17    float cut = cell.getCut();
18    for(auto p= startPtr; p<endPtr; ++p)
19    {
20        nLeft += *p < cut;
21    }
22
23    out[cellPtrOffset] = nLeft;
24 }
```

Combining the data is straight forward, as in *Combine()* *vout* and *vout2* contain the counts for a given number of particles. As the master is interested in a global count across all the particles contained within all local storages, we simply sum each element of *vout* together with *vout2*.

Listing 5: Part of the Count Left Combine() method

```

1 for(auto i=0; i<nCounts; ++i)
2     out[i] += out2[i];
```

5.5.3 Count

Counting the total number of particles for each cell is trivial with the *cellToRange* map, as we can directly know the result when subtracting the begin index from the end index of each cell. The *Combine()* is equal to the one shown for the Count Left Service.

Listing 6: Part of Count Service() method

```

1 for (int cellPtrOffset=0; cellPtrOffset<nCells; ++cellPtrOffset) {
2     auto cell = static_cast<Cell>*(in + cellPtrOffset);
3     out[cellPtrOffset] =
4     lcl->cellToRangeMap(cell.id,1) -
5     lcl->cellToRangeMap(cell.id,0);
6 }
```

5.5.4 Partition

The Partitioning Service is a direct implementation of the hoare partition as seen in listing 7. Again the service receives a buffer with a given number of cells, which are then iterated over using the outer for loop. Each cell from the buffer is then being casted to the proper struct (line 2). Using the id of the cell, we can then access the start and end indices of the corresponding particles stored in local storage (lines 4-5). Lines 7-37 search for pairs located on the wrong side of the cut and then swap them using a simple *swap()* helper function. *i* has to be initialized using *beginInd* - 1 as the *i++* operation happens before a read on the array (line 15). For the same reason *j* is initialized to *endInd* and not *endInd* - 1. Then the while loop terminates, *i* points to the first index in the array where the value is greater than the cut, therefore we can update the *cellToRangeMap* with *i*, as for the left child cell the start remains equal but the end is now *i* and for the right child cell its vice versa.

Listing 7: Partition Service

```

1 for (int cellPtrOffset=0; cellPtrOffset<nCells; ++cellPtrOffset){
2     auto cell = static_cast<Cell>*(in + cellPtrOffset);
3
4     int beginInd = pst->lcl->cellToRangeMap(cell.id, 0);
5     int endInd = pst->lcl->cellToRangeMap(cell.id, 1);
6
7     int i = beginInd-1, j = endInd;
8     float cut = cell.getCut();
9
10    while(true)
11    {
12        do
13        {
14            i++;
15        } while(lcl->particles(i, cell.cutAxis)<cut and i<=endInd);
16
17        do
18        {
19            j--;
20        } while(lcl->particles(j, cell.cutAxis)>cut and j>=beginInd);
21
22        if(i >= j) {
23            break;
24        }
25
26        swap(lcl->particles, i, j);
27    }
28
29    swap(lcl->particles, i, endInd -1);
30
31    lcl->cellToRangeMap(cell.getLeftChildId(), 0) =
32    lcl->cellToRangeMap(cell.id, 0);
33    lcl->cellToRangeMap(cell.getLeftChildId(), 1) = i;
34
35    lcl->cellToRangeMap(cell.getRightChildId(), 0) = i;
36    lcl->cellToRangeMap(cell.getRightChildId(), 1) =
37    lcl->cellToRangeMap(cell.id, 1);
38
39 }

```

5.5.5 Make Axis

As mentioned in section 2.1 the relevant axis to search for the cut position differs for each cell. Therefore we cannot simply iterate over a single array, we need to iterate over the relevant axis instead. In order to simplify the process we introduce a service which copies iterates over all cells and copies the slice of coordinates which are encompassed in its volume and lie on its cut axis to a temporary array.

5.6 Parallel Schedule

In the context of parallelization, we define the number of processors as np . Initially we assume that each processor has a random unique subset of all the particles stored

in its local memory, this has two reasons. For one we run into memory limitations quickly when trying to load all the particles onto a single node, furthermore memory bandwidths limitations can be multiplied by the number of processors and are thus a lot higher.

In the running example this might look as follows:

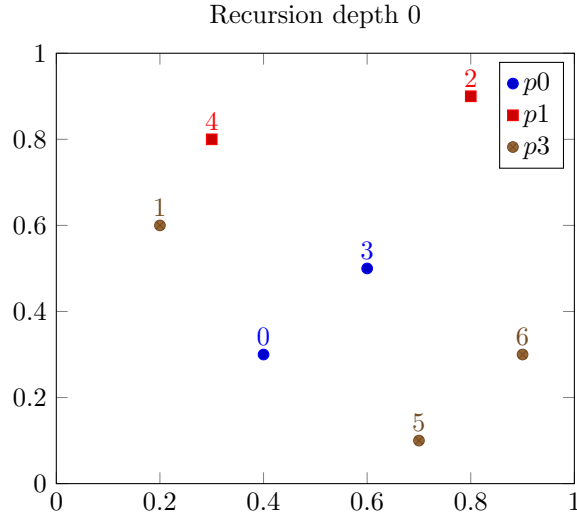


Figure 18: Example particles distributed randomly across 3 nodes

The Master schedules all other processors but also performs tasks by itself. The exact schedule is illustrated in figure 19 where a service which is dispatched from the master and executed on all processors is depicted as a horizontal rectangle. Computations which are only performed by a single processor are represented by a vertical rectangle.

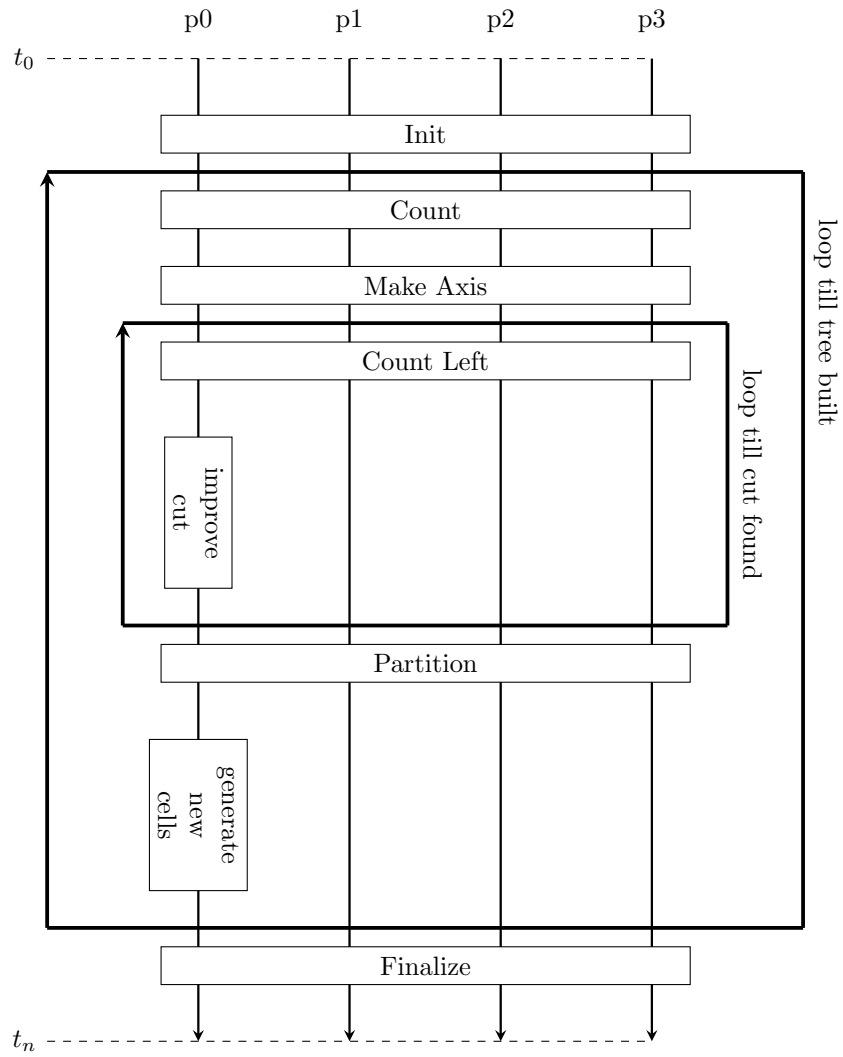


Figure 19: Parallelized ORB CPU version

6 GPU Implementation

CUDA is a parallel computing platform and programming model invented by *NVIDIA*. It allows software developers and programmers to use a CUDA-enabled graphics processing unit (GPU) for general purpose processing. The CUDA platform is designed to work with programming languages such as C, C++, and Fortran and gives direct access to the GPU's virtual instruction set and parallel computational elements.

6.1 Relevant CUDA Concepts

CUDA gives us the ability to launch kernels, which are written in C with some additional CUDA specific syntax. These kernel can be run from the CPU, commonly referred to as the host, and are then executed on the GPU, also known as the device. The device and the host have a separate memory and as seen in section 4.1 the transfer I_{CG} rates between the GPU and CPU are generally slower, under performing I_C and I_G speeds. Therefore one of key the challenges when rewriting CPU code to GPU code is to limit data transfers between the GPU and CPU as much as possible. Furthermore we have divide our problems into subproblems, where each subproblem is then executed on a block. CUDA cannot give any guarantees considering the order of execution of these blocks, therefore we have to fundamentally rethink algorithms when porting them from CPU to GPU code, where some algorithms are more or less fitted. Generally we only think about problems which are applied on large data arrays, the less connection between the individual results exists, the easier it is to implement an GPU version of the computation. Furthermore we can assume that problems with less branching work better on the GPU due to various reasons, which we will explain in detail later.

We can access the location of a thread within a block using *threadIdx.x* as well as the block ID among all other blocks associated with the same kernel *blockIdx.x*. Finally the number of blocks are stored in the variable *gridDim.x*.

Building a tree is therefore a rather challenging problem as there are many computations which influence another and building a tree involves a lot of branching, at least when it is done in a conventional manner.

Listing 8: Calling a CUDA Kernel

```
1 add<<<<
2     nBlocks ,
3     nThreads ,
4     nThreads * bytesPerThread ,
5     stream
6 >>>>(
7     param1 ,
8     param2
9 );
```

A kernel is executed exactly once for every kernel, where a block consists of many kernels. The number of kernels per block and the total number of blocks can be defined by the user as seen in the kernel call syntax in figure 8 on line 2 and 3.

6.1.1 Warps

Each consecutive grouping of 32 threads form a warp. All threads within a warp are executed in parallel, given that there is no warp divergence. Because warps are

executed in parallel, there is no need for synchronization between the threads of the same warp. A warp divergence can occur, if there is a control statements, where two or more threads from the same warp execute a different code. This leads to a decrease in performance and should be avoided whenever possible.

6.1.2 Memory

Each block has access to shared memory register, where the capacity of this register can be defined at runtime by the host as seen in figure 8 on line 4. The maximum shared memory size depends depending on the architecture ranges between 48 and 96kB. [14] There are also upper limits for the number of threads per block which is usually between 512 and 1024. [14] With a shared memory size of 48kb and 1024 threads per block we get $\frac{48000B}{1024} = 46B$ per thread, which is roughly equivalent to 11 32 bit precision numbers per thread. However as we usually do not max out the thread counts per block, there is more shared memory space available per thread.

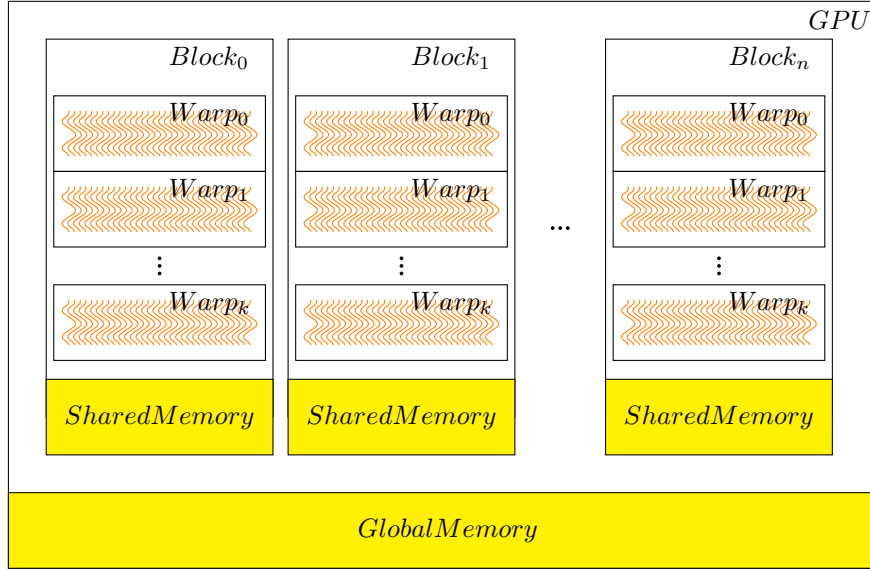
Data in the shared memory exists only as long as the kernel exists and it cannot be accessed across different blocks. Race conditions do apply to shared memory, but threads in a block can be synchronized, therefore it is not safe to have writes from multiple threads on the same shared memory address. There exists a CUDA command `__syncthreads()` which allows us to synchronize all threads within a block, enabling us to control the execution order of individual statements.

Additionally to the shared memory, there exists global memory. Global memory is persistent even after a block or a kernel has finished and is the only way to send data between the device and the host. Since the data is not deleted after a kernel has finished, we can and should reuse the data whenever possible for other kernels, reducing the memory transfer overhead. Global memory is significantly slower than shared memory and it is best practice to copy data from global memory to shared memory before performing actual computations on it, thus reducing the total number of memory accesses on global memory. After the computations have finished, the results can be written back from shared to global memory.

6.1.3 Memory Access Patterns

Global memory is fetched from memory in 32 Byte packages, which translates to $\frac{32}{4} = 8$ floats. A simple aligned sequential memory access pattern, where each thread reads a single float from global memory, results in 4 transaction per warp (32 threads). Unaligned memory access can have minor performance penalties because more memory banks have to be loaded per warp.

For shared memory each bank consists of 32 bits, which translates to a 32 bit precision number. If 32 unique threads access 32 different banks, then its considered an ideal access pattern, because the banks can be loaded in parallel. If however multiple threads access the same memory, a bank conflict is introduced and the access speed is greatly reduced.



6.1.4 Asynchronous Operations

There are two types of engines that can be used to execute kernels in CUDA streams: copy engines and kernel engines. Copy engines are used to copy data between host and device memory, and between different types of memory on the device. Kernel engines are used to execute CUDA kernels. The number of individual engines depends on the actual hardware, lower end GPUs usually have a single kernel and a copy engine, more advanced architectures can have more than one copy or kernel engine.

A CUDA stream is a sequence of commands that are executed in order on a CUDA device. Streams can be used to improve the runtime of a CUDA program by overlapping the execution of different kernels. For example, a copy kernel can be executed in one stream while a compute kernel is executed in another stream. This overlap can lead to a significant performance improvement.

6.1.5 Synchronization

Whilst threads within a block can be synchronized, this is not the case for the blocks launched from a kernel. The only real synchronization technique, is to wait for a kernel, meaning all blocks associated with this kernel, to finish and execute a consecutive task using another kernel.

Since the GPU is limited in its computing and also memory capabilities, the number of blocks which are run in parallel are limited, therefore CUDA cannot give us any guarantee, whether a set of blocks are run serially or in parallel. Some of the blocks might be run in parallel, meanwhile others are run serially. This very constraint, forces

the programmer to rethink algorithms and think of ways, how individual subproblems can be run somewhat independently of each other.

There exists however a way to communicate between blocks in a safe way. Atomic operations can be performed on global memory without any race conditions.

6.2 Streams

As described in section 6.1.4, we can leverage streams to improve the runtime a CUDA accelerated application. The way we make use of streams is simple, each thread performs all its computations using on a unique stream. Since leveraging streams only improves the runtime in low numbers because the actual copy and kernel engines are not very numerous, this variant is very simple and effective. We have observed that using more streams on the individual threads does not improve the runtime, quiet contrary it can have a negative impact.

6.3 Memory Management

We use several strategies in order to improve the performance of memory allocation:

1. **Pinned Memory** As the GPU cannot access the default paged memory directly, when copying memory from the host to the device, the memory must first be copied to pinned memory. This means if we use pinned memory directly with *cudaMallocHost()*, data transfers can be around twice as fast. In our implementation we use pinned memory for all data, which need to be either sent from the host to device or vice versa.
2. **Reuse Memory** We can avoid allocation and freeing commands altogether and instead reuse previously allocated memory whenever possible. In our implementation we allocate the memory in the very beginning. Since we have a fixed number of particles as well as an upper limit for the number of cells memory can be allocated when calling the initialize service. The concept is applied to both CPU and GPU memory.

6.4 GPU Accelerated Count Left

We have described a CPU implementation of the Count Left Service. In this section we explain how we can port this problem to the GPU. In essence we can use a general reduction as a basis version, and adapt it to sum up elements which fulfill a certain condition, where the condition is being smaller than a given value. The code is based on the reduction as its explained in a Webinar from *NVIDIA*. [15]

6.4.1 Schedule

The entire schedule of the ORB is depicted in figure 20. In a first step we call the initialize service, where the necessary data is allocated on all devices and the particles are loaded. Furthermore the initial SPTDS is constructed, which essentially only consists of the root node, encompassing the entire domain and all the particles. Next we enter the main loop of ORB, it iterates until it has constructed a SPTDS with the desired size. Within the loop the count service is called, computing the number of cells for each cell summed over all the threads in the system. To prepare the data for the GPU transfer we make the temporary array using the make axis service, which is

then sent to the GPU. We can now start with the root finding process, where we count the particles left of the initial cut position using the Count Left Service. Depending on the outcome, we then improve the cut position and repeat until a nearly perfect cut is found for all the leaf cells of the current SPTDS. We now generate two new child cells using the computed cut positions for all leaf cells of the current SPTDS and enrich the current SPTDS with the newly generated cells. Finally the particles array is partitioned accordingly. The tree building loop is then repeated or if the desired size of the SPTDS is reached, we exit the loop and call the finalize service to free the allocated memory.

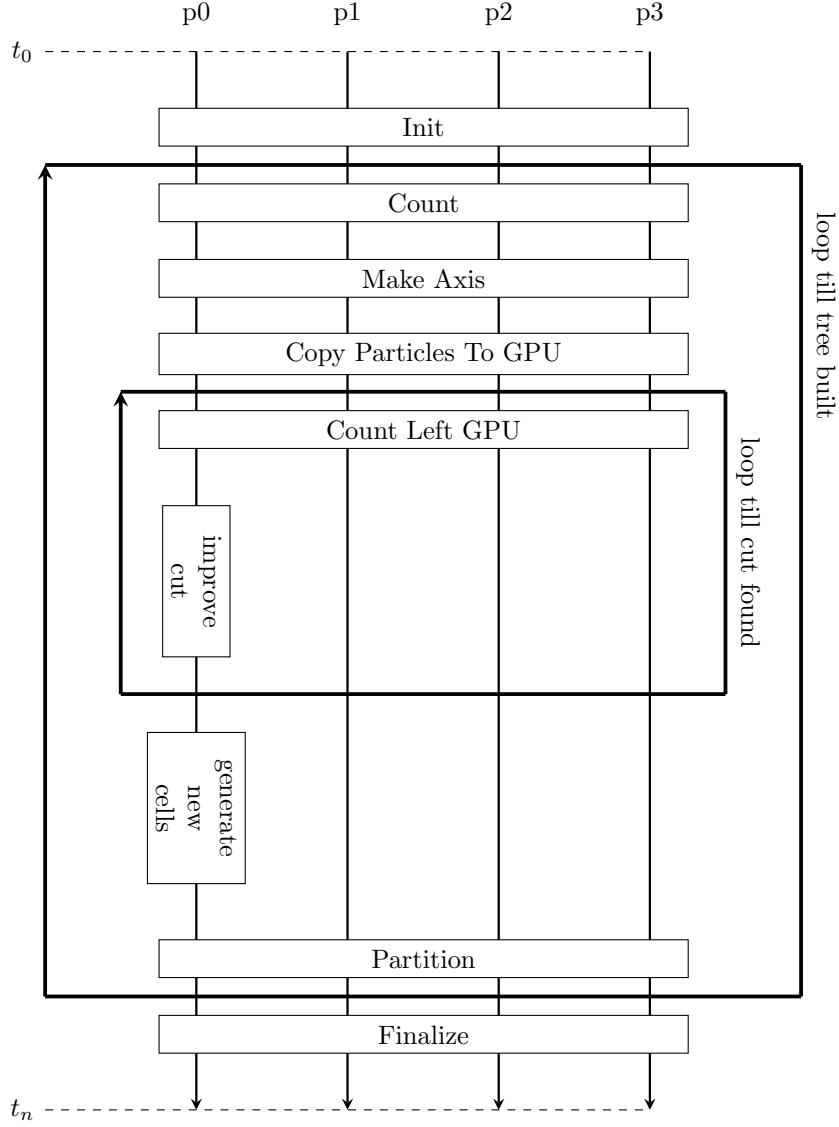


Figure 20: Parallelized ORB GPU version

6.4.2 Service

The GPU version of the Count Left Service invokes the reduction kernel exactly once for every leaf cell of the current SPTDS. Meaning on level 0 of the tree there is only one kernel per thread being executed, when proceeding to further levels, the number of kernels is equivalent to 2^{depth} . The number of blocks can be determined by the number of particles in this cell n_{cell} divided by the number of threads per block which

is generally set to 256 and finally we divide this further by the number of elements per thread.

Since we can pass all the relevant information regarding a cell as input parameters to the cell, there is no need to copy any additional data from the CPU to the GPU. We only need to copy back the results from the GPU to the GPU after the kernel has finished.

6.4.3 Kernel Code

The actual kernel code is depicted in listing 9. The input parameters of the reduction kernel (line 13-16) are in order: the input data, which is in sequential order the particles, the output array, where the results of the kernel is stored, the position of the cut plane and finally the total number of elements contained within the cell. Note that since we make use of the Make Axis Service, there is no need to know to which axis the orthogonal cut plane is searched for.

In a first step on lines 18-20 each thread ID *tid* is associated with an index in the global memory of input data. Furthermore on line 22-23 the shared data is initialized and set to zero. Its size is equivalent to the block size, the total number of thread per block. Meaning we can store a single float per thread. On lines 25-28 the thread iterates over all elements which it is associated to and copies the results of a compare to shared memory. This allows us to work exclusively with shared memory instead of global memory, reducing the amount of costly memory accesses.

There exists an optimum in terms of number of operation performed by a single thread. In order to adapt the number of ops dynamically, the while loop ensures that each thread iterates over a certain number of elements as seen in lines 25-28. The input parameter *n* defines the total number of particles. Therefore if we reduce the total number of blocks by a factor of *r*, the while loop will iterate over *r* elements. Inside the while loop on line 27 the index *i* is incremented by *gridSize* in each iteration. The grid size is equivalent to the total number of threads associated with the kernel. Therefore when launching the kernel, the number of blocks can be divided by *r* which then results in *r* iterations performed in the while loop. In our case we set *r* to 16 as more or less items per thread will result in a lower performance.

The rest of the kernel code performs the actual reduction, where in each step two elements from the shared memory are summed together. Given a block size of 256, after line 37 exactly 256 sums are stored in shared memory. In a next step as seen on line 37 (lines 31-36 are not relevant as the size is smaller than 512) the left half threads in the block add their value together with the values stored in a position offset by 128. This step reduces the total number of elements to be considered to 128, and the step can be repeated as seen on lines 42 to 46.

The reduction is optimized by using a template, which indicates the total number of thread per block, which is equivalent to the blockSize. Since the template is evaluated at compile time, all the if statements taking blockSize as a parameter in figure do not cause any performance loss. In fact, because the reduction loops are unrolled (31-52) the the performance of the kernel is increased.

Each time there is a control statement involving non template arguments, a branch divergence is introduced. On the GPU these divergences become problematic bad when they are introduced inside a warp. Meaning if any of the 32 within a warp execute a different code from the other threads, the warp cannot function optimally anymore. For this reason, the *warpReduce* method as seen on line 2-9 is introduced. This method is only executed when there are 64 elements in shared memory which are

still to be summed, which requires only 32 threads as every thread in each iteration sums two elements together. Inside the warp reduce method, there are more unrolled loops continuing the same patterns as described above. However in this case there are no control statements that could introduce a warp divergence, which means there are computations performed which are not actually necessary, however as the computations inside a warp are executed in parallel anyways, this comes with no performance penalties. Inside a warp there is no need for synchronization as all threads are executed in parallel.

Listing 9: Conditional Reduction in CUDA

```

1 template <unsigned int blockSize>
2 extern __device__ void warpReduce(
3     volatile int *sdata, unsigned int tid) {
4     if (blockSize >= 64) sdata[tid] += sdata[tid + 32];
5     if (blockSize >= 32) sdata[tid] += sdata[tid + 16];
6     if (blockSize >= 16) sdata[tid] += sdata[tid + 8];
7     if (blockSize >= 8) sdata[tid] += sdata[tid + 4];
8     if (blockSize >= 4) sdata[tid] += sdata[tid + 2];
9     if (blockSize >= 2) sdata[tid] += sdata[tid + 1];
10 }
11
12 template <unsigned int blockSize>
13 extern __global__ void reduce(
14     float *g_idata,
15     unsigned int *g_odata,
16     float cut,
17     int n) {
18     unsigned int tid = threadIdx.x;
19     unsigned int i = blockIdx.x*(blockSize) + threadIdx.x;
20     unsigned int gridSize = blockSize*gridDim.x;
21
22     __shared__ int sdata[blockSize];
23     sdata[tid] = 0;
24
25     while (i < n) {
26         sdata[tid] += (g_idata[i] <= cut);
27         i += gridSize;
28     }
29     __syncthreads();
30
31     if (blockSize >= 512) {
32         if (tid < 256) {
33             sdata[tid] += sdata[tid + 256];
34         }
35         __syncthreads();
36     }
37     if (blockSize >= 256) {
38         if (tid < 128) {
39             sdata[tid] += sdata[tid + 128];
40         } __syncthreads();
41     }
42     if (blockSize >= 128) {
43         if (tid < 64) {
44             sdata[tid] += sdata[tid + 64];
45         } __syncthreads();
46     }

```

```

47  if (tid < 32) {
48      warpReduce<blockSize>(sdata, tid);
49  }
50  if (tid == 0) {
51      g_odata[blockIdx.x] = sdata[0];
52  }
53 }

```

We execute the reduction kernel once for each cell, where we distribute the particles within the cell among a set of blocks. This makes the implementation straight forward but as we increase the number of cells, we make more calls to the reduction kernel. This is problematic because each initialization of a kernel comes with some overhead, degrading the performance gradually with the tree traversal, as the overhead to computation ratio becomes more and more unfavorable.

6.5 Improved GPU Accelerated Count Left

To solve the degrading performance problem we make some changes to the kernel and the overall schedule. The main idea is to prepare the necessary data for all cells in advance, effectively reducing the overall number of necessary kernel calls to $\lceil \log_2 d \rceil$ instead of $2 \times d$. Each block is provided with information concerning the start and end index of the cells particles, the cut position as well as an index of the cell.

We will explain the improved version using a simple example: Let us consider *cell*₂ with a volume that encompasses particles in the range 0 - 10240. *cell*₃ encompasses particles in the range 10240 - 20480. We have a blocksize of 256 threads and the elements per thread are 16, thus each block processes $256 * 16 = 4096$ elements. In this case we reserve three blocks for each cell, where the first block processes elements 0 - 4096 and the second one 4096 - 10240. By assigning more work to the second block, we make sure that there are no underworked threads, as measurements have shown a slight increase in elements per thread is better than a decrease. Furthermore we can avoid blocks where not even all the threads are occupied, because less than 256 elements are processed on the block. Consequently the third and fourth block respectively process particles 10240 to 14446 and 14446 to 20480.

6.5.1 Schedule

Only a single change has to be made to the schedule: We introduce another service which is called GPU Copy Cell Service, which does as the name says, copying the cell information from the CPU to the GPU. Since only the cut positions change while we are iterating over the inner loop, this service can be called outside of it. The Improved GPU Count Left Service copies the cut data each time it is invoked. Other than that, the schedule remains equal.

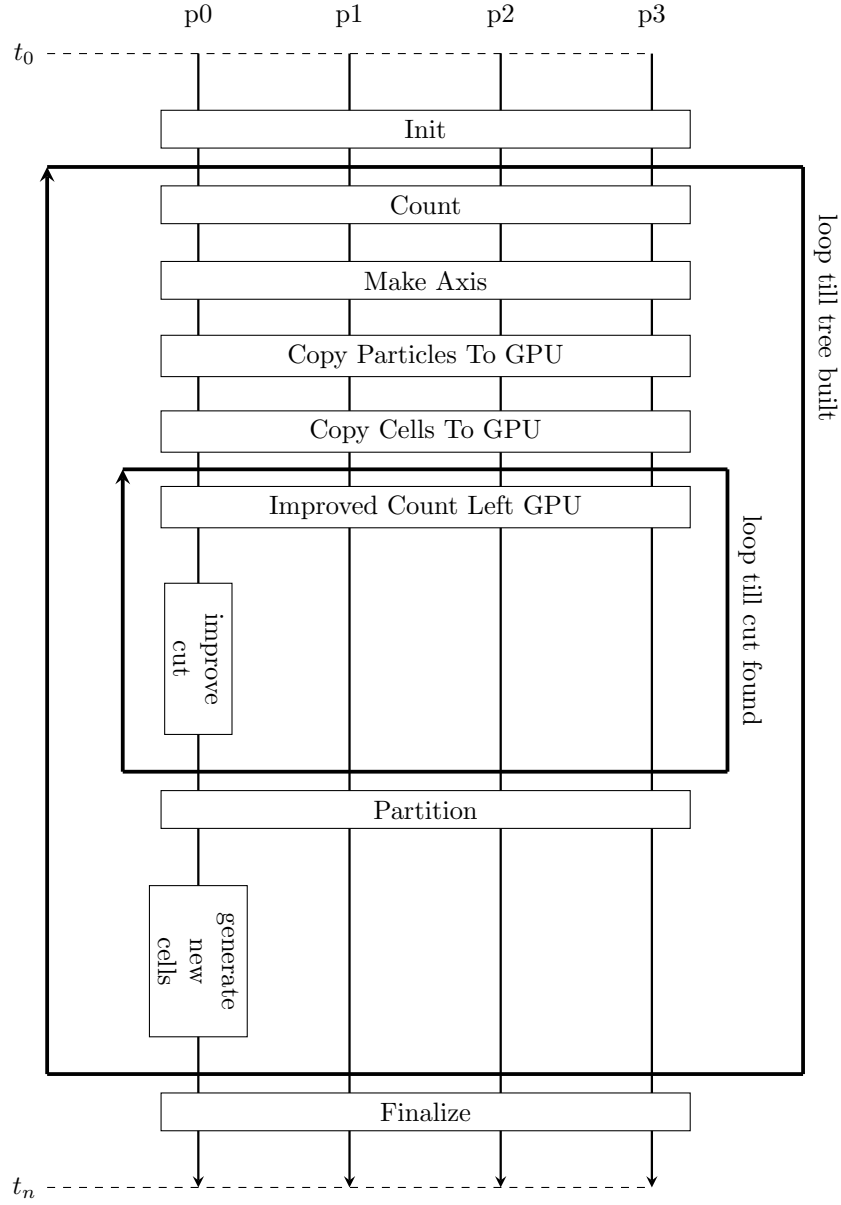


Figure 21: Parallelized ORB GPU version 2

6.5.2 Kernel Code

The kernel code as seen in listing 10 largely remains the same. New data pointers are passed along as a parameter where g_{begin} and g_{end} mark the begin and end of

the important particle array slice, furthermore g_{cuts} is an array of cell cuts for each block. As we prepare the data in a way, that each block only needs to operate over the particles contained within a single cell, the data can be read using the `blockIdx.x` (lines 20 - 22) which is a CUDA variable and provides a sequential indexing of blocks associated with the kernel. Everything else, including the subroutine *warpReduce* remain equal.

Listing 10: Kernel Optimized GPU Count Left

```

1  template <unsigned int blockSize>
2  extern __global__ void reduce(
3  float * g_idata ,
4  unsigned int * g_begins ,
5  unsigned int * g_ends ,
6  float * g_cuts ,
7  unsigned int * g_odata) {
8      __shared__ unsigned int s_data[blockSize];
9
10     unsigned int tid = threadIdx.x;
11     const unsigned int begin = g_begins[blockIdx.x];
12     const unsigned int end = g_ends[blockIdx.x];
13     const float cut = g_cuts[blockIdx.x];
14
15     unsigned int i = begin + tid;
16     s_data[tid] = 0;
17
18     // unaligned coalesced g memory access
19     while (i < end) {
20         s_data[tid] += (g_idata[i] <= cut);
21         i += blockSize;
22     }
23     __syncthreads();
24
25     if (blockSize >= 512) {
26         if (tid < 256) {
27             s_data[tid] += s_data[tid + 256];
28         }
29         __syncthreads();
30     }
31     if (blockSize >= 256) {
32         if (tid < 128) {
33             s_data[tid] += s_data[tid + 128];
34         } __syncthreads();
35     }
36     if (blockSize >= 128) {
37         if (tid < 64) {
38             s_data[tid] += s_data[tid + 64];
39         } __syncthreads();
40     }
41     if (tid < 32) {
42         warpReduce<blockSize>(s_data, tid);
43     }
44     if (tid == 0) {
45         g_odata[blockIdx.x] = s_data[0];
46     }
47 }

```

Figure 22: Reduction in CUDA

6.6 GPU Accelerated Partitioning

We have successfully implemented a version with a GPU accelerated Count Left. Since the particles are still partitioned on the CPU, the data needs to be copied to GPU each time we expand the SPTDS by another . As we use the Make Axis service to compress the axes into a single array, we only need to copy N elements, but we need to do so $\lceil \log_2 d \rceil$ times. If we partition the particles on the GPU, the data needs to be copied only once, but in this case all three coordinates have to be copied, thus $N \times 3$. Nevertheless this corresponds a fraction of the total data transfer: $\frac{3}{\lceil \log_2 d \rceil}$. In figure 23 its observable how after d reaches 2^3 the number of bytes to be transferred are less

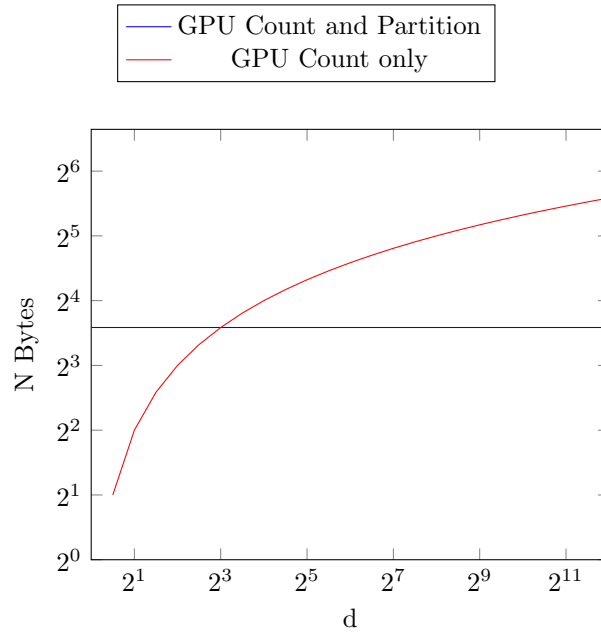


Figure 23: Data transferred with and without GPU partitioning

A GPU partitioning algorithm was first proposed by Cederman [16] and later refined by Manca [17] where atomic operations are leveraged to improve the runtime of the algorithm. Our version is built around the concepts introduced by Manca.

Early results have shown partitioning the data on the GPU is not really faster than the CPU version and the number of particles which can be stored on the GPU memory is reduced by $\frac{5}{3}$. This is caused by required additional temporary arrays for memory swaps on the GPU. Furthermore there was not enough time to fully optimize and track down all issues with the code, therefore we do not describe the full kernel code and instead limit ourselves to lay the ground work for subsequent work.

6.6.1 Schedule

In the adapted schedule ?? we can move the Copy Particles To GPU service out of the main loop, other than that the schedule remains the same.

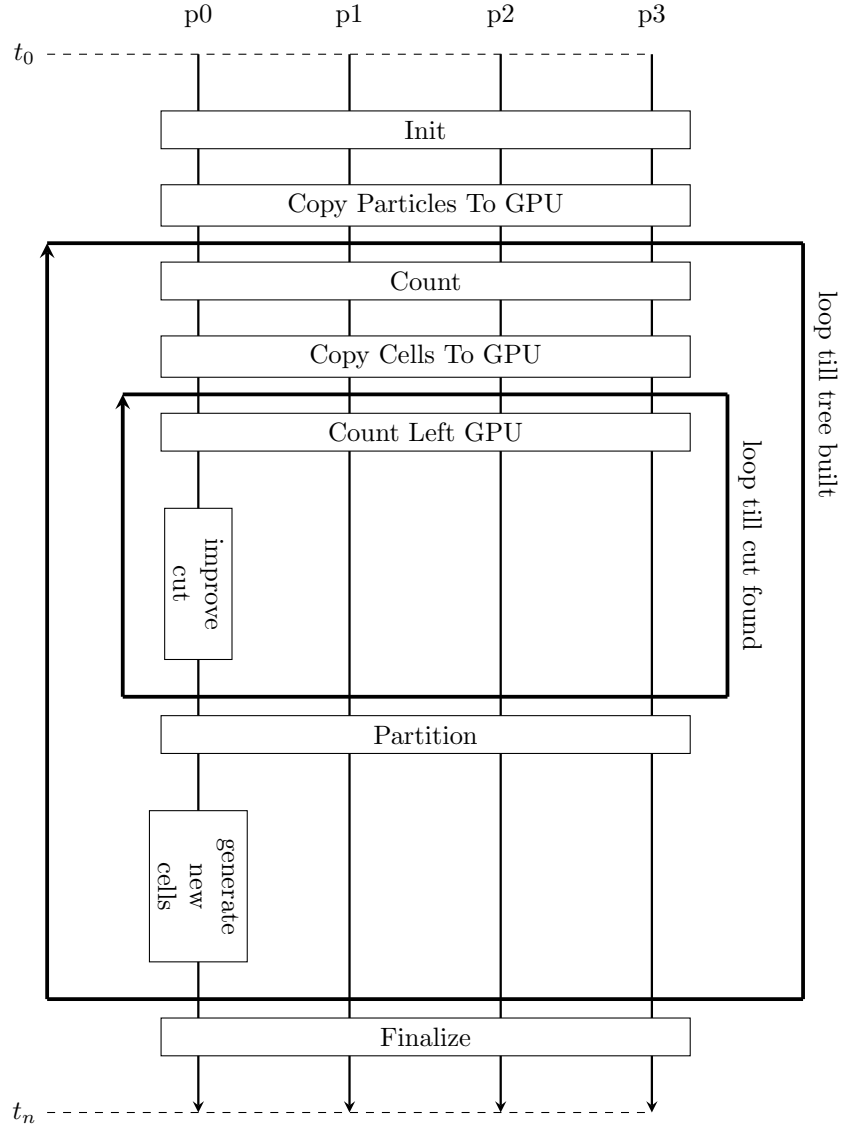


Figure 24: Parallelized ORB with GPU counting and GPU partitioning

6.6.2 Partitioning a Single Cell

In a first step we perform a block wise partition, meaning there is no synchronization between blocks at this stage. Each thread loads the particle coordinate of the matching axis into the shared memory register S . We then make two copies of the data stored in the shared memory, the first copy A compares each value with the cut and stores 0 if the value is smaller and 1 otherwise. The second copy B stores the inverted result. We now perform a parallel prefix sum, or scan, on both copies. We can use an algorithm as it is described in GPU Gems [18]. On B we add the total number of elements smaller than the cut to each element. In theory the block wise partition could already be applied by storing each element with index i larger than the cut at A_i and larger elements at B_i . However our goal is to partition all elements across multiple blocks. To do so we introduce a global variable called leq and g where the total number of elements smaller equal and bigger than the cut are stored. Each time a block has computed both scans, the number of elements designated to the block smaller and larger than the cut are known. Thus a single thread of each block performs an atomic add on leq and g accordingly. $AtomicAdd()$ returns the old value of the variable. Therefore we can use and store it in two designated shared memory register s_{leq} and s_g . At this stage every thread in a block agrees on a value for s_{leq} and s_g the values can be treated as starting indices to store the block wise partitioned results in the output array. In this case however we store an element larger than the cut at position i in location $A_i + s_{leq}$ and otherwise at $B_i = s_g$.

In our case the input array are coordinate values of particles and the output array is a temporary array of size N . The in an output array cannot be the same, meaning the partition cannot be performed in place, as doing so would result in elements that are overwritten before they are read by the thread assigned to it. This cannot be avoided as the relationship between the range of input data treated by a kernel and the output is not deterministic. Its dictated by the order, by which the kernel reach the $atomicAdd()$ statement.

As of now only a single axis is partitioned, but the same needs to be done for the other two axes. One approach would be to partition all the other arrays in the same kernel, as the permutations can simply be reapplied. However in this case we need three additional temporary arrays, reducing the number of particles which can be stored on the GPU by $\frac{1}{2}$. A more memory efficient, but slightly slower variant is to allocate another array on the GPU where the actual permutation indices are stored. Therefore together with the temporary array N is only reduced by $\frac{3}{5}$. After the partition kernel has finished, we can copy back the coordinates from the temporary array to the original location. Then the second axis is partitioned using the precomputed permutations and the results are again copied back to the proper location. Finally after the step is repeated for the last axis, the partition is finished.

6.6.3 Partitioning Multiple cells

We apply the same method as introduced in section 6.5 to ensure only a single kernel has to be launched for all the cells.

In a first step each thread loads the particle coordinate of the matching axis into the shared memory register.

shared memory of each block in a primary step. This can be achieved by copying all data from global to shared memory and then performing two scans over the elements. Initiall we scan the elements smaller than the cut, secondly we perform a scan of the

elements which are larger. The results of scanning the elements smaller than the cut can directly be considered as index positions in the block sized subproblem. Equally each of the scan results of the elements larger than the cut, plus the total

The problem with partitioning the array on the GPU is that it reduces the number of particles which can be stored in GPU memory by $\frac{5}{3}$. Lets us consider we want to partition all particles ranging from 0 to N along the x axis with a cut plane positioned at 0. In a first step we compute

The partitioning algorithm is a widely used concept and there are various implementations used in different variants of the quick-sort algorithm. The overall design of the implementation is guided by **TODO: Insert ref** and the scan is taken from GPU Gems **TODO: Insert ref**. The main and probably only disadvantage of partitioning the particle array on the GPU, is that it cannot be done in place, thus we need another temporary array allocated on the GPU memory which further restricts the total number of particles. Furthermore we need to store the permutations, because we need to apply them to the other axes as well. We could as well simply perform all permutations from the same kernel at the same time, but this would require 3 temporary arrays, thus reducing the total number of particles by a factor of 2. The solution with the partition array and a single temporary array only results in a reduction of the particles by a factor of $\frac{5}{3}$.

7 Performance Analysis of ORB

All measurements were performed on Piz Daint where the performance of system can be read from table 1. The supercomputer uses a slurm workload manager. Simple job script can be executed using *sbatch*. In our case we used a meta script which has the function of generating the actual job scripts for the desired parameters.

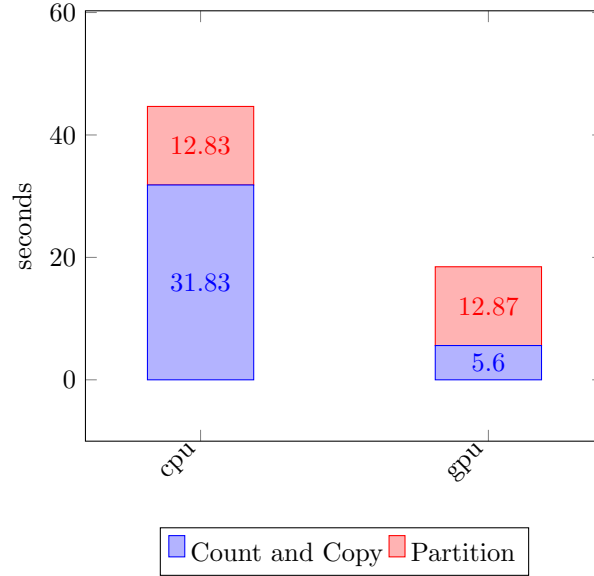


Figure 25: Execution times of different strategies

In figure 25 the final measurements are plotted. The total runtime of the ORB is around 43 seconds when performed only on the CPU, for the GPU accelerated version it comes down to 18 seconds. However a large part of the execution time can be attributed to the Partitioning, where its overhead remains equal for both versions. Even the CPU Partitioning can be improved a lot and the GPU version might increase the performance even further. If we subtract the time lost for the partitioning, the GPU version is about 5.6 times faster, which is very close to the 5.25 which were expected in the theoretical model. The total runtime is not equal to the one expected, however the theoretical model works with upper bound runtimes for the bisection method, in the method terminates at around 20 iterations. When scaling the theoretical model with this, the results look a lot closer.

7.1 Scaling

To make sure the GPU version scales good we have evaluated the total runtimes for

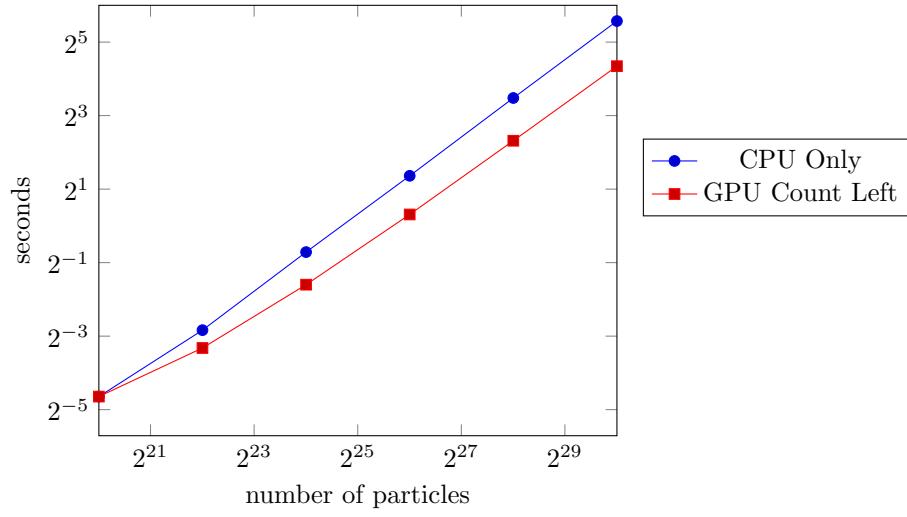


Figure 26: Execution times of different strategies

The total execution time is more than halved. Note that this also includes the unoptimized partition part. If we repeat the plots but without the partition

7.2 Comparison to Theoretical Model

8 Conclusion

References

- [1] Joachim Gerhard Stadel. *Cosmological N-body simulations and their analysis*. PhD thesis, 2001.
- [2] ICS Software, August 2022. [Online; accessed 1. Aug. 2022].
- [3] TOP500, August 2022. [Online; accessed 1. Aug. 2022].
- [4] Günther Voglsam. Real-time ray tracing on the gpu - ray tracing using cuda and kd-trees. Master’s thesis, Institute of Computer Graphics and Algorithms, Vienna University of Technology, April 2013.
- [5] Thrust, May 2022. [Online; accessed 28. Jul. 2022].
- [6] CUB, May 2022. [Online; accessed 28. Jul. 2022].
- [7] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms*. MIT press, 2022.
- [8] Piz Daint & Piz Dora, April 2022. [Online; accessed 27. Apr. 2022].
- [9] ORNL Launches Summit Supercomputer | ORNL, April 2022. [Online; accessed 27. Apr. 2022].
- [10] Intel® Xeon® Processor E5-2690 , February 2022. [Online; accessed 27. Apr. 2022].
- [11] AMD EPYC™ 7742 | AMD, April 2022. [Online; accessed 27. Apr. 2022].
- [12] NVIDIA Tesla P100: The Most Advanced Data Center Accelerator, April 2022. [Online; accessed 27. Apr. 2022].
- [13] Todd Veldhuizen. Blitz++ User’s Guide , May 2022. [Online; accessed 11. May 2022].
- [14] CUDA C++ Programming Guide, July 2022. [Online; accessed 2. Aug. 2022].
- [15] Mark Harris. Optimizing parallel reduction in cuda, August 2022. [Online; accessed 1. Aug. 2022].
- [16] Daniel Cederman and Philippas Tsigas. GPU-Quicksort: A practical Quicksort algorithm for graphics processors. *ACM J. Exp. Algorithmics*, 14:4–1, January 2010.
- [17] Emanuele Manca, Andrea Manconi, Alessandro Orro, Giuliano Armano, and Luciano Milanesi. CUDA-quicksort: an improved GPU-based implementation of quicksort. *Concurrency Computat.: Pract. Exper.*, 28(1):21–43, January 2016.
- [18] Chapter 39. Parallel Prefix Sum (Scan) with CUDA, August 2022. [Online; accessed 3. Aug. 2022].

List of Figures

1	Uniform random distribution of 3D coordinates in cube domain projected onto a 2d plane	9
2	Example distribution with $N = 7$	13
3	Example particles with ORB at recursion depth 0	13

4	Tree with ORB at recursion depth 0	13
5	Example particles with ORB at recursion depth 1	14
6	Tree with ORB at recursion depth 1	14
7	Example particles with ORB at recursion depth 2	15
8	Tree with ORB at recursion depth 2	15
9	Illustration of the universal computing model	21
10	Roofline Model for Piz Daint CPU	24
11	Roofline Model for Piz Daint GPU	25
12	Roofline Model for Summit	25
13	Roofline Model for Summit GPU	26
14	Roofline Model for Alps	26
15	Execution times of different strategies	30
16	Tree as heap	32
17	Tree with $d = 3$	33
18	Example particles distributed randomly across 3 nodes	38
19	Parallelized ORB CPU version	39
20	Parallelized ORB GPU version	45
21	Parallelized ORB GPU version 2	49
22	Reduction in CUDA	51
23	Data transferred with and without GPU partitioning	52
24	Parallelized ORB with GPU counting and GPU partitioning	53
25	Execution times of different strategies	56
26	Execution times of different strategies	57

List of Tables

1	Datapoints of Supercomputers	22
---	--	----